

P-Values for Classification – Computational Aspects and Asymptotics

Inauguraldissertation
der Philosophisch-naturwissenschaftlichen Fakultät
der Universität Bern
und
der Fakultät für Mathematik und Informatik
der Georg-August-Universität Göttingen

vorgelegt von
Niki Roger Zumbrunnen
von Aeschi bei Spiez

Leiter der Arbeit:
Prof. Dr. L. Dümbgen
Institut für mathematische Statistik und Versicherungslehre
der Universität Bern
und
Prof. Dr. A. Munk
Institut für mathematische Stochastik
der Georg-August-Universität Göttingen

Originaldokument gespeichert auf dem Webserver der Universitätsbibliothek Bern



Dieses Werk ist unter einem
Creative Commons Namensnennung-Keine kommerzielle Nutzung-Keine Bearbeitung 2.5
Schweiz Lizenzvertrag lizenziert. Um die Lizenz anzusehen, gehen Sie bitte zu
<http://creativecommons.org/licenses/by-nc-nd/2.5/ch/> oder schicken Sie einen Brief
an
Creative Commons, 171 Second Street, Suite 300, San Francisco, California 94105, USA.

Urheberrechtlicher Hinweis

Dieses Dokument steht unter einer Lizenz der Creative Commons
Namensnennung-Keine kommerzielle Nutzung-Keine Bearbeitung 2.5
Schweiz.

<http://creativecommons.org/licenses/by-nc-nd/2.5/ch/>

Sie dürfen



dieses Werk vervielfältigen, verbreiten und öffentlich zugänglich machen.

Zu den folgenden Bedingungen:



Namensnennung. Sie müssen den Namen des Autors/Rechteinhabers in der von ihm festgelegten Weise nennen (wodurch aber nicht der Eindruck entstehen darf, Sie oder die Nutzung des Werkes durch Sie würden entlohnt).



Keine kommerzielle Nutzung. Dieses Werk darf nicht für kommerzielle Zwecke verwendet werden.



Keine Bearbeitung. Dieses Werk darf nicht bearbeitet oder in anderer Weise verändert werden.

Im Falle einer Verbreitung müssen Sie anderen die Lizenzbedingungen, unter welche dieses Werk fällt, mitteilen.

Jede der vorgenannten Bedingungen kann aufgehoben werden, sofern Sie die Einwilligung des Rechteinhabers dazu erhalten.

Diese Lizenz lässt die Urheberpersönlichkeitsrechte nach Schweizer Recht unberührt.

Eine ausführliche Fassung des Lizenzvertrags befindet sich unter
<http://creativecommons.org/licenses/by-nc-nd/2.5/ch/legalcode.de>

P-Values for Classification – Computational Aspects and Asymptotics

Inauguraldissertation
der Philosophisch-naturwissenschaftlichen Fakultät
der Universität Bern
und
der Fakultät für Mathematik und Informatik
der Georg-August-Universität Göttingen

vorgelegt von
Niki Roger Zumbrunnen

von Aeschi bei Spiez

Leiter der Arbeit:
Prof. Dr. L. Dümbgen
Institut für mathematische Statistik und Versicherungslehre
der Universität Bern

und
Prof. Dr. A. Munk
Institut für mathematische Stochastik
der Georg-August-Universität Göttingen

Von der Philosophisch-naturwissenschaftlichen Fakultät angenommen.

Bern, 05. 03. 2014

Der Dekan:
Prof. Dr. S. Decurtins

Abstract

P-Values for Classification. Let (\mathbf{X}, Y) be a random variable consisting of an observed feature vector \mathbf{X} and an unobserved class label $Y = 1, 2, \dots, L$ with unknown joint distribution. In addition, let \mathcal{D} be a training data set consisting of n completely observed independent copies of (\mathbf{X}, Y) .

First, we consider a point predictor for Y , namely the standard linear classifier for two classes. But we do not assume Gaussian distributions. In this setting we provide a central limit theorem for missclassification rates and cross-validated estimators thereof.

Point predictors do not provide information about confidence. To get such information, we construct for each $b = 1, 2, \dots, L$ a p-value $\pi_b(\mathbf{X}, \mathcal{D})$ for the null hypothesis that $Y = b$, treating Y temporarily as a fixed parameter, i.e. we construct a prediction region for Y with a certain confidence. In particular, we consider p-values based on the plug-in statistic for the standard model with two classes and prove a central limit theorem for inclusion probabilities and cross-validated estimators thereof.

In addition, we discuss data-driven choices of tuning parameters for p-values based on multicategory logistic regression, where we use regularization terms to deal with high-dimensional feature vectors \mathbf{X} .

Randomized P-Values. Randomized tests are a familiar concept from mathematical statistics. The goal is to obtain tests with exact prescribed significance level even in settings with test statistics having discrete distributions. We discuss the related concept of randomized p-values. One benefit is that p-values obtained from different independent test statistics can be combined more easily. Since in applications non-randomized tests and p-values are needed, we review and modify a method of Meinshausen et al. (2009).

A major example is the analysis of several independent contingency tables with small cell counts. We propose various ways of combining corresponding randomized p-values. We also illustrate the benefits of the final de-randomized test.

Acknowledgements

First and foremost, I would like to thank the main supervisor, Lutz Dümbgen, for his excellent support throughout my studies. His always-open office door, his positive attitude and his patience made it a great pleasure to work with him.

A special thanks goes to the co-supervisor, Axel Munk, for his valuable inputs and for giving me the opportunity of spending a month at the University of Göttingen. I am also grateful to Prof. Dr. Regina Liu for reviewing my thesis.

Many thanks also to Jasmin Wandel, Dominic and Heike Schuhmacher, Benjamin Baumgartner, David Ginsbourger, Chris Kopp, Bernhard Freiermuth and all other members of the IMSV for helpful and interesting discussions, table football and other enjoyable breaks.

Finally, I would like to thank my family and Olivia Jutzi for their support and encouragement.

This work was supported by the research group FOR916 of the Swiss National Science Foundation (SNF) and the Deutsche Forschungsgemeinschaft (DFG).

Contents

Overview	1
1. Classifiers and P-Values	3
1.1. Classification	3
1.1.1. Optimal Classifiers in the Ideal Case	3
1.1.2. Classification Using Training Data	5
1.1.3. Estimation of Missclassification Rates	9
1.2. From Classifiers to P-Values	9
1.3. Optimal P-Values as Benchmark	11
1.4. P-Values via Permutation Tests	15
1.5. Estimation of Separability	17
1.6. Asymptotic Properties	18
1.7. Implementation in <code>pvc</code>	20
1.7.1. Shortcut	21
1.7.2. Data Example ‘buerk’	21
1.7.3. Main Functions	23
1.8. Technical Details for Penalized Multicategory Logistic Regression	27
1.8.1. The Log-Likelihood-Function	27
1.8.2. Regularizations	31
1.8.3. Strict Convexity and Coercivity	35
1.8.4. Some Comments on the Implementation in <code>pvc</code>	36
2. Choice of Tuning Parameters	39
2.1. Stability	39
2.1.1. Subsampling	40
2.1.2. Extended Golden Section Search	40
2.2. Dimension Reduction	41
2.3. Numerical Examples	42
2.3.1. Simulated Data	42
2.3.2. Real Data	45
3. Central Limit Theorems	47
3.1. Half-Spaces	48
3.1.1. Root- n -Consistency	48
3.1.2. Empirical Processes	54

3.2. Asymptotics of Estimators for Location and Scatter	55
3.3. A Central Limit Theorem for Missclassification Rates	59
3.4. A Central Limit Theorem for Inclusion Probabilities	70
4. Randomized and De-Randomized P-Values	83
4.1. De-Randomization	83
4.2. Combining Independent P-Values	86
4.3. Application to Multiple Contingency Tables	88
4.3.1. Two-by-Two Tables	88
4.3.2. K-by-L Tables	92
A. Classical Results	97
A.1. Lindeberg-Feller Central Limit Theorem	97
A.2. Neyman-Pearson Lemma	97
References	99
Index	101
List of Symbols	103

Overview

We start in Section 1.1 with a short introduction to classification. In the remaining part of Chapter 1 we present the p-values for classification introduced by Dümmbgen et al. (2008). First we assume that the joint distribution of (\mathbf{X}, Y) is known. In this setting optimal p-values are available. If the joint distribution is unknown, we use training data to compute nonparametric p-values based on permutation tests. We review asymptotic results of Dümmbgen et al. (2008) and Zumbunnen (2009). Finally, we comment on technical details and the implementation of the p-values in the R-package `pvclass`.

Some of the test statistics we use depend on a tuning parameter such as the k in the nearest neighbor method or the penalty parameter τ in the logistic regression. In Chapter 2 we propose a data-driven way to choose such parameters. In addition, we comment on computational issues.

The theoretic main result is given in Chapter 3. First we consider linear discriminant analysis with two classes. But we do not assume Gaussian distributions. To estimate the covariance matrix we use either the standard estimator or more robust M -estimators. In these two settings we present central limit theorems for missclassification rates and cross-validated estimators thereof. This result implies in particular that the estimators are root- n -consistent.

Next we consider p-values based on the plug-in statistic for the standard model with two classes. But we relax the assumption of Gaussian distributions to elliptical symmetry. The corresponding conditional inclusion probabilities are of interest to judge the separability of the two classes. However, these theoretic quantities are typically unknown. Therefore we use cross-validation to estimate them. Dümmbgen et al. (2008) proved that these estimators are consistent. We take a closer look at the inclusion probabilities and the corresponding estimators and describe their asymptotic distribution. In particular, we derive a central limit theorem, which implies that the estimators are root- n -consistent. Moreover, it enables us to construct confidence intervals for the inclusion probabilities.

For the computation of the p-values, we add the new observation temporarily to a certain class. But it may be an outlier with respect to the distribution of this class. Therefore it is reasonable to use a robust M -estimator for the covariance matrix. Our asymptotic results are valid for the standard estimator as well as for the M -estimators.

Chapter 4 is not directly related to the main part of this thesis. In this chapter we discuss the concept of randomized p-values. Since in applications non-randomized tests and p-values are needed, we review and modify a method of Meinshausen et al. (2009). We propose various ways of combining corresponding randomized p-values. We also illustrate the benefits of the final de-randomized test for several independent contingency tables with small cell counts.

1. Classifiers and P-Values

In this chapter we first give a short introduction to classification. For a more detailed introduction and further references we refer to McLachlan (1992). Then we introduce nonparametric p-values for classification as they are given in Dümmbgen et al. (2008) and comment on the implementation in the R-package `pvc`lass.

This chapter is mainly based on Zumbrennen (2009) and Dümmbgen (2011).

1.1. Classification

Let (\mathbf{X}, Y) be a pair of random variables, where

$$Y \in \mathcal{Y} := \{1, \dots, L\}, \quad L \geq 2$$

is the *class label* of an observation, which is described by the *feature vector* \mathbf{X} with values in the *feature space* \mathcal{X} .

Classifying (\mathbf{X}, Y) means that only \mathbf{X} is observed and Y has to be predicted via \mathbf{X} .

1.1.1. Optimal Classifiers in the Ideal Case

Suppose that the joint distribution of (\mathbf{X}, Y) is known, i.e. we know the *prior probabilities*

$$w_\theta := \mathbb{P}(Y = \theta)$$

and the conditional distributions

$$P_\theta := \mathcal{L}(\mathbf{X} \mid Y = \theta)$$

for all $\theta \in \mathcal{Y}$. Further let M be a measure on \mathcal{X} and let each conditional distribution P_θ be described by a density function f_θ with respect to M , i.e.

$$P_\theta(B) = \int_B f_\theta(\mathbf{x}) M(d\mathbf{x}),$$

for measurable sets $B \subset \mathcal{X}$.

1. Classifiers and P-Values

In the simplest case a *classifier* is a point predictor $\hat{Y}(\mathbf{X}): \mathcal{X} \rightarrow \mathcal{Y}$ for Y . To find an optimal classifier we need a quality criterion, for example the risk of missclassification

$$R(\hat{Y}) := \mathbb{P}(\hat{Y}(\mathbf{X}) \neq Y).$$

The following proposition characterizes *optimal classifiers* $\hat{Y}^*(\mathbf{X})$, in the sense that they minimize the risk R :

Lemma 1.1. *Let $\hat{Y}(\mathbf{X}): \mathcal{X} \rightarrow \mathcal{Y}$ be a classifier. Then*

$$R(\hat{Y}) \geq 1 - \int \max_{\theta \in \mathcal{Y}} w_{\theta} f_{\theta}(\mathbf{x}) M(d\mathbf{x}),$$

with equality if and only if

$$\hat{Y}(\mathbf{x}) \in \arg \max_{\theta \in \mathcal{Y}} w_{\theta} f_{\theta}(\mathbf{x}) \quad \text{for } M\text{-almost all } \mathbf{x} \in \mathcal{X}. \quad (1.1)$$

PROOF. With $\mathbb{1}(A)$ denoting the indicator function for the set A ,

$$\begin{aligned} R(\hat{Y}) &:= \mathbb{P}(\hat{Y}(\mathbf{X}) \neq Y) \\ &= 1 - \sum_{\theta \in \mathcal{Y}} w_{\theta} \mathbb{P}(\hat{Y}(\mathbf{X}) = \theta \mid Y = \theta) \\ &= 1 - \sum_{\theta \in \mathcal{Y}} w_{\theta} \int f_{\theta}(\mathbf{x}) \mathbb{1}\{\hat{Y}(\mathbf{x}) = \theta\} M(d\mathbf{x}) \\ &= 1 - \int w_{\hat{Y}(\mathbf{x})} f_{\hat{Y}(\mathbf{x})}(\mathbf{x}) M(d\mathbf{x}) \\ &\geq 1 - \int \max_{\theta \in \mathcal{Y}} w_{\theta} f_{\theta}(\mathbf{x}) M(d\mathbf{x}). \end{aligned}$$

The inequality is obviously an equality if and only if (1.1) is satisfied. \square

Standard Gaussian Model

Let $P_{\theta} = \mathcal{N}_d(\boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma})$ with mean vectors $\boldsymbol{\mu}_{\theta} \in \mathbb{R}^d$ and a common symmetric, positive definite covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$.

We may write the Gaussian densities as

$$f_{\theta}(\mathbf{x}) = c \exp(-D_{\boldsymbol{\Sigma}}^2(\mathbf{x}, \boldsymbol{\mu}_{\theta})/2)$$

with $c := (2\pi)^{-d/2} \det(\boldsymbol{\Sigma})^{-1/2}$ and the *Mahalanobis distance*

$$D_{\boldsymbol{\Sigma}}(\mathbf{x}, \mathbf{y}) := \sqrt{(\mathbf{x} - \mathbf{y})^{\top} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \mathbf{y})} = \|\boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \mathbf{y})\|$$

between $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^d$ (with respect to Σ). Here and throughout $\|\cdot\|$ denotes the Euclidean norm for vectors or the Frobenius norm for matrices, respectively.

Therefore the optimal classifier can be characterized by

$$\hat{Y}^*(\mathbf{x}) \in \arg \min_{\theta \in \mathcal{Y}} (D_{\Sigma}^2(\mathbf{x}, \mu_{\theta}) - 2 \log w_{\theta}). \quad (1.2)$$

Suppose that $L = 2$ and $\mu_1 \neq \mu_2$. Then

$$D_{\Sigma}^2(\mathbf{x}, \mu_1) - 2 \log w_1 \begin{cases} > \\ = \\ < \end{cases} D_{\Sigma}^2(\mathbf{x}, \mu_2) - 2 \log w_2$$

if and only if

$$(\mathbf{x} - \mu_{1,2})^{\top} \Sigma^{-1} (\mu_2 - \mu_1) \begin{cases} > \\ = \\ < \end{cases} \log(w_1/w_2),$$

with $\mu_{\theta,b} := (\mu_{\theta} + \mu_b)/2$ for $\theta, b \in \mathcal{Y}$. The sets of all feature vectors assigned uniquely to class 1 or 2, respectively, are separated by a hyperplane orthogonal to $\Sigma^{-1}(\mu_2 - \mu_1)$.

1.1.2. Classification Using Training Data

The joint distribution of (\mathbf{X}, Y) is typically unknown. In this case we estimate the prior probabilities w_{θ} and the densities f_{θ} by adequate estimators $\hat{w}_{\theta}(\mathcal{D})$ and $\hat{f}_{\theta}(\cdot, \mathcal{D})$, respectively. Then we choose a classifier

$$\hat{Y}(\mathbf{x}, \mathcal{D}) \in \arg \max_{\theta \in \mathcal{Y}} \hat{w}_{\theta} \hat{f}_{\theta}(\mathbf{x}).$$

To estimate w_{θ} and f_{θ} we use *training data* \mathcal{D} , consisting of pairs (\mathbf{X}_i, Y_i) , for $i = 1, \dots, n$, whereas the \mathbf{X}_i as well as the Y_i are known. We consider the (\mathbf{X}_i, Y_i) as random variables with the same distribution as (\mathbf{X}, Y) , and assume that the $n+1$ data pairs $(\mathbf{X}, Y), (\mathbf{X}_1, Y_1), (\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_n, Y_n)$ are stochastically independent. Let

$$\mathcal{G}_{\theta} := \{i \leq n : Y_i = \theta\} \quad \text{and} \quad N_{\theta} := \#\mathcal{G}_{\theta}.$$

Then an estimator for w_{θ} is given by

$$\hat{w}_{\theta} := \frac{N_{\theta}}{n}.$$

Linear Discriminant Analysis

In the standard model with $P_\theta = \mathcal{N}_d(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma})$ we replace the unknown mean vectors $\boldsymbol{\mu}_\theta$ and covariance matrix $\boldsymbol{\Sigma}$ in (1.2) with corresponding estimators and get the *standard linear classifier*

$$\hat{Y}^*(\mathbf{x}) \in \arg \min_{\theta \in \mathcal{Y}} (D_{\hat{\boldsymbol{\Sigma}}}^2(\mathbf{x}, \hat{\boldsymbol{\mu}}_\theta) - 2 \log \hat{w}_\theta). \quad (1.3)$$

The standard estimators in this model are given by

$$\begin{aligned} \hat{\boldsymbol{\mu}}_\theta &:= \frac{1}{N_\theta} \sum_{i \in \mathcal{G}_\theta} \mathbf{X}_i, \\ \hat{\boldsymbol{\Sigma}} &:= \frac{1}{n-L} \sum_{\theta \in \mathcal{Y}} \sum_{i \in \mathcal{G}_\theta} (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_\theta)(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_\theta)^\top. \end{aligned}$$

The assumption of Gaussian distributions could be relaxed, e.g. to elliptically symmetric distributions.

Elliptically Symmetric Distributions. The random vector $\mathbf{Z} \in \mathbb{R}^d$ has a *spherically symmetric* distribution (with respect to $\mathbf{0}$) if for any orthonormal matrix $\mathbf{B} \in \mathbb{R}^{d \times d}$, $\mathbf{B}\mathbf{Z}$ has the same distribution as \mathbf{Z} . The distribution of a random vector $\mathbf{X} \in \mathbb{R}^d$ is *elliptically symmetric* with center $\boldsymbol{\mu} \in \mathbb{R}^d$ and scatter matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$, if $\boldsymbol{\Sigma}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$ has a spherically symmetric distribution.

For a spherically symmetric random vector $\mathbf{Z} \in \mathbb{R}^d$ with $\mathbb{P}(\mathbf{Z} = \mathbf{0}) = 0$ and any unit vector $\mathbf{v} \in \mathbb{R}^d$

$$\mathcal{L}(\mathbf{v}^\top \mathbf{Z}) = \mathcal{L}(Z_1), \quad (1.4)$$

where Z_1 is the first component of \mathbf{Z} . For the proof of this claim and more details on elliptically symmetric distributions we refer to Muirhead (1982).

Robust M-Estimators. The standard estimator for $\boldsymbol{\Sigma}$ is sensitive to outliers. As an alternative, we consider the more robust M -estimators $\hat{\boldsymbol{\Sigma}}_M$ and $\hat{\boldsymbol{\Sigma}}_{sym}$. The first M -estimator, $\hat{\boldsymbol{\Sigma}}_M$, is the maximum likelihood estimator in the model where $P_\theta = \mathcal{N}_d(\boldsymbol{\mu}_\theta, c_\theta \boldsymbol{\Sigma})$ with $c_\theta > 0$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ symmetric and positive definite with $\det(\boldsymbol{\Sigma}) = 1$. For the calculations we use that $\hat{\boldsymbol{\Sigma}}_M$ is the solution of the fixed point equation

$$\boldsymbol{\Sigma} = d \sum_{\theta=1}^L \frac{N_\theta}{n} \frac{\mathbf{M}_\theta}{\text{trace}(\boldsymbol{\Sigma}^{-1} \mathbf{M}_\theta)}$$

with $\mathbf{M}_\theta := \sum_{i \in \mathcal{G}_\theta} (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_\theta)(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_\theta)^\top$.

The second M -estimator, $\hat{\Sigma}_{sym}$, is a generalization for more than one group of the symmetrized version of Tyler's M -estimator, as it is defined in Dümbgen (1998). We assume that the observations \mathbf{X}_i are pairwise different within groups. Then $\hat{\Sigma}_{sym}$ is the solution of the fixed point equation

$$\Sigma = \frac{d}{c} \sum_{\theta=1}^L \frac{1}{N_\theta} \sum_{\substack{i,j \in \mathcal{G}_\theta \\ i < j}} \frac{(\mathbf{X}_i - \mathbf{X}_j)(\mathbf{X}_i - \mathbf{X}_j)^\top}{(\mathbf{X}_i - \mathbf{X}_j)^\top \Sigma^{-1} (\mathbf{X}_i - \mathbf{X}_j)}$$

with $c := \sum_{\theta=1}^L (N_\theta - 1)/2 = (n - L)/2$.

k Nearest Neighbors

The nearest-neighbor method is a nonparametric alternative to estimate f_θ . It has the advantage, that only few assumptions about the distributions P_θ are required. Suppose that (\mathcal{X}, d) is a separable metric space and consider the closed balls

$$B(\mathbf{x}, r) := \{\mathbf{y} \in \mathcal{X} : d(\mathbf{x}, \mathbf{y}) \leq r\}$$

and the open balls

$$U(\mathbf{x}, r) := \{\mathbf{y} \in \mathcal{X} : d(\mathbf{x}, \mathbf{y}) < r\}$$

for all $\mathbf{x} \in \mathcal{X}$ and $r \geq 0$. Assume that

$$M(B(\mathbf{x}, r)) < \infty \quad \text{for all } \mathbf{x} \in \mathcal{X} \text{ and } r \geq 0$$

and that f_θ is continuous for all $\theta \in \mathcal{Y}$.

Lemma 1.2. *Let \mathcal{X}_0 be the support of $\mathcal{L}(\mathbf{X})$, i.e.*

$$\mathcal{X}_0 := \{\mathbf{x} \in \mathcal{X} : \mathbb{P}(\mathbf{X} \in B(\mathbf{x}, r)) > 0 \text{ for all } r > 0\}.$$

Then \mathcal{X}_0 is the smallest closed set with $\mathbb{P}(\mathbf{X} \in \mathcal{X}_0^c) = 0$ and

$$f_\theta(\mathbf{x}) = \lim_{r \downarrow 0} \frac{P_\theta(B(\mathbf{x}, r))}{M(B(\mathbf{x}, r))} \quad \text{for all } \theta \in \mathcal{Y}, \mathbf{x} \in \mathcal{X}_0. \quad (1.5)$$

To estimate P_θ , we use the *empirical measure* of the points $\mathbf{X}_i, i \in \mathcal{G}_\theta$,

$$\hat{P}_\theta(B) := \frac{\#\{i \in \mathcal{G}_\theta : \mathbf{X}_i \in B\}}{N_\theta} \quad \text{for } B \subset \mathcal{X}.$$

Now define for fixed integer $k \leq n$ and any $\mathbf{x} \in \mathcal{X}$

$$\hat{r}_{k,n}(\mathbf{x}) = \hat{r}_{k,n}(\mathbf{x}, \mathcal{D}) := \min \{r \geq 0 : \#\{i \leq n : \mathbf{X}_i \in B(\mathbf{x}, r)\} \geq k\}$$

1. Classifiers and P-Values

such that $B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x}))$ is the smallest ball centered at \mathbf{x} , which covers at least k training vectors \mathbf{X}_i . These are the k nearest neighbors of \mathbf{x} . Next we define

$$\hat{f}_\theta(\mathbf{x}) := \frac{\hat{P}_\theta(B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x})))}{M(B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x})))}.$$

Since $M(B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x})))$ is the same for all classes $\theta \in \mathcal{Y}$, we end up with an estimator of the form

$$\hat{Y}_k(\mathbf{x}) \in \arg \max_{\theta \in \mathcal{Y}} \hat{w}_\theta \hat{P}_\theta(B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x}))).$$

For $\hat{w}_\theta = N_\theta/n$, this can be written as

$$\hat{Y}_k(\mathbf{x}) \in \arg \max_{\theta \in \mathcal{Y}} \#\{i \in \mathcal{G}_\theta : \mathbf{X}_i \in B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x}))\}.$$

This means we use majority vote among the k nearest neighbors to classify \mathbf{X} .

PROOF OF LEMMA 1.2. First we show that \mathcal{X}_0 is closed. Let $\mathbf{x} \in \mathcal{X} \setminus \mathcal{X}_0$. Then there is a $r > 0$ such that $\mathbb{P}(\mathbf{X} \in U(\mathbf{x}, r)) = 0$. Now let $\mathbf{y} \in U(\mathbf{x}, r)$. Then $U(\mathbf{y}, \tilde{r}) \subset U(\mathbf{x}, r)$ with $\tilde{r} = r - d(\mathbf{x}, \mathbf{y}) > 0$. But this implies that $\mathbb{P}(\mathbf{X} \in U(\mathbf{y}, \tilde{r})) = 0$, and thus $\mathbf{y} \notin \mathcal{X}_0$. Therefore $U(\mathbf{x}, r) \cap \mathcal{X}_0 = \emptyset$. Since the choice of $\mathbf{x} \in \mathcal{X} \setminus \mathcal{X}_0$ was arbitrary, this implies that \mathcal{X}_0 is closed.

Let \mathcal{X}_* be a countable and dense subset of \mathcal{X} . For each $\mathbf{x} \in \mathcal{X}_0^c$ there exists an $r > 0$ such that $\mathbb{P}(\mathbf{X} \in B(\mathbf{x}, r)) = 0$. Then there is a $\mathbf{x}_* \in \mathcal{X}_*$ such that $d(\mathbf{x}, \mathbf{x}_*) < r/2$, so $\mathbf{x} \in B(\mathbf{x}_*, r_*) \subset B(\mathbf{x}, r)$ for any $r_* \in (r/2, r - d(\mathbf{x}, \mathbf{x}_*)) \cap \mathbb{Q}$. This construction shows that

$$\mathcal{X}_0^c \subset \bigcup \{B(\mathbf{x}_*, r_*) : \mathbf{x}_* \in \mathcal{X}_*, 0 < r_* \in \mathbb{Q}, \mathbb{P}(\mathbf{X} \in B(\mathbf{X}_*, r_*)) = 0\}.$$

Since the latter union is countable, it has measure zero as well.

Now suppose that there is a closed set $\mathcal{X}'_0 \subsetneq \mathcal{X}_0$ satisfying $\mathbb{P}(\mathbf{X} \in \mathcal{X}'_0^c) = 0$. Then for any $\mathbf{x} \in \mathcal{X}_0 \setminus \mathcal{X}'_0$ there is a $r > 0$ such that $B(\mathbf{x}, r) \cap \mathcal{X}'_0 = \emptyset$ and so $\mathbb{P}(\mathbf{X} \in B(\mathbf{x}, r)) = 0$, which contradicts the definition of \mathcal{X}_0 . Thus \mathcal{X}_0 is the smallest closed set with $\mathbb{P}(\mathbf{X} \in \mathcal{X}_0^c) = 0$.

Next we show that (1.5) holds. Since

$$P_\theta(B(\mathbf{x}, r)) = \int_{B(\mathbf{x}, r)} f_\theta(\mathbf{y}) M(d\mathbf{y}) \begin{cases} \leq \sup_{\mathbf{y} \in B(\mathbf{x}, r)} f_\theta(\mathbf{y}) M(B(\mathbf{x}, r)) \\ \geq \inf_{\mathbf{y} \in B(\mathbf{x}, r)} f_\theta(\mathbf{y}) M(B(\mathbf{x}, r)), \end{cases}$$

it follows from the continuity of f that

$$\lim_{r \downarrow 0} \frac{P_\theta(B(\mathbf{x}, r))}{M(B(\mathbf{x}, r))} = f_\theta(\mathbf{x}).$$

□

Weighted Nearest Neighbors

In the previous paragraph the k -nearest neighbors of the observation \mathbf{X} were all weighted equally. However, it would be reasonable to assign larger weights to training observations which are closer to \mathbf{X} . Now we first order the training data according to their distance to \mathbf{X} and then we assign descending weights to them. Let

$$W_n(1) \geq W_n(2) \geq \dots \geq W_n(n) \geq 0$$

be weights and

$$\hat{R}(\mathbf{x}, \mathbf{X}_i) := \#\{j \leq n : d(\mathbf{x}, \mathbf{X}_j) \leq d(\mathbf{x}, \mathbf{X}_i)\}$$

the rank of training observation \mathbf{X}_i according to its distance to \mathbf{x} . If two or more training observations are at equal distance to \mathbf{x} we could average the weights or use randomization. However, we assume that

$$M\{\mathbf{y} \in \mathcal{X} : d(\mathbf{x}, \mathbf{y}) = r\} = 0 \quad \text{for all } \mathbf{x} \in \mathcal{X}, r \geq 0 \quad (1.6)$$

to avoid this problem. We define the *weighted nearest neighbor classifier* by

$$\hat{Y}_{wnn}(\mathbf{x}) \in \arg \max_{\theta \in \mathcal{Y}} \sum_{i \in \mathcal{G}_\theta} W_n(\hat{R}(\mathbf{x}, \mathbf{X}_i)).$$

1.1.3. Estimation of Missclassification Rates

To judge a certain classifier one could estimate the missclassification rates

$$R_\theta = \mathbb{P}(\hat{Y}(\mathbf{X}, \mathcal{D}) \neq Y \mid Y = \theta, \mathcal{D})$$

using cross-validation, i.e. with the estimator

$$\hat{R}_\theta = \frac{\#\{i \in \mathcal{G}_\theta : \hat{Y}(\mathbf{X}_i, \mathcal{D}_i) \neq Y_i\}}{N_\theta},$$

where \mathcal{D}_i denotes the training data without observation (\mathbf{X}_i, Y_i) , and compare them with the missclassification rates of the optimal classifier.

1.2. From Classifiers to P-Values

A drawback of point estimators is the lack of information about confidence. To get such information we could use a Bayesian approach and calculate the *posterior distribution* of Y given \mathbf{X} , i.e. the *posterior weights*

$$w_\theta(\mathbf{X}) := \mathbb{P}(Y = \theta \mid \mathbf{X}).$$

1. Classifiers and P-Values

By Lemma 1.1, a classifier \hat{Y}^* satisfying

$$\hat{Y}^*(\mathbf{x}) \in \arg \max_{\theta \in \mathcal{Y}} w_{\theta}(\mathbf{x})$$

is optimal in the sense, that it minimizes the risk $R(\hat{Y}) = \mathbb{P}(\hat{Y}(\mathbf{X}) \neq Y)$. Thus we can now compute the conditional risk

$$\mathbb{P}(\hat{Y}^*(\mathbf{X}) \neq Y \mid \mathbf{X} = \mathbf{x}) = 1 - \max_{\theta \in \mathcal{Y}} w_{\theta}(\mathbf{x}),$$

which gives us information about confidence of \hat{Y}^* .

However a drawback of the posterior probabilities is, that the posterior weights $w_{\theta}(\mathbf{X})$ depend sensitively on the prior weights w_{θ} , i.e. small changes in w_{θ} can lead to totally different $w_{\theta}(\mathbf{X})$, which we illustrate in Example 1.1. In addition, classes with small prior weights w_{θ} tend to be ignored by the classifier \hat{Y}^* , and so the class-dependent risk $\mathbb{P}(\hat{Y}^*(\mathbf{X}) \neq Y \mid Y = \theta)$ may be rather large for some classes θ . Moreover, in some studies the class labels are not random, but predetermined by the study design. For example in a case-control study, one recruits a certain number of diseased individuals and a certain number of healthy individuals. Furthermore, in medical studies the prior probabilities can change over time or differ geographically, while the distributions P_{θ} are reasonably assumed to be universal. Another problem arises, if the future observation (\mathbf{X}, Y) belongs to a so far unknown class $\theta \notin \mathcal{Y}$.

In the daily routine one often uses a process of elimination to classify objects. In our context, this means that we exclude certain classes $\theta \in \mathcal{Y}$ and finally give a set of plausible candidates for Y . In other words we treat Y temporarily as an unknown fixed parameter and compute for each class $\theta \in \mathcal{Y}$ a p-value $\pi_{\theta}(\mathbf{X})$ or $\pi_{\theta}(\mathbf{X}, \mathcal{D})$ for the null hypothesis that $Y = \theta$. In the ideal case, where the joint distribution of (\mathbf{X}, Y) is known, this means $\pi_{\theta}: \mathcal{X} \rightarrow [0, 1]$ satisfies

$$\mathbb{P}(\pi_{\theta}(\mathbf{X}) \leq \alpha \mid Y = \theta) \leq \alpha \quad \text{for all } \alpha \in (0, 1). \quad (1.7)$$

Given such p-values π_{θ} , the set

$$\hat{\mathcal{Y}}_{\alpha}(\mathbf{X}) := \{\theta \in \mathcal{Y}: \pi_{\theta}(\mathbf{X}) > \alpha\}$$

is a $(1 - \alpha)$ -prediction region for Y , i.e.

$$\mathbb{P}(Y \in \hat{\mathcal{Y}}_{\alpha}(\mathbf{X}) \mid Y = \theta) \geq 1 - \alpha \quad \text{for any } \theta \in \mathcal{Y}, \alpha \in (0, 1).$$

Thus we can exclude the classes $\theta \notin \hat{\mathcal{Y}}_{\alpha}(\mathbf{X})$ with confidence $1 - \alpha$. If there is only one $\theta \in \hat{\mathcal{Y}}_{\alpha}(\mathbf{X})$, we have classified \mathbf{X} uniquely with confidence $1 - \alpha$.

Since we compute p-values for multiple null hypotheses, one could expect that we get a multiple testing problem. However, the observation \mathbf{X} belongs to only one class and therefore at most one null hypothesis is true.

In the realistic case, where the joint distribution of (\mathbf{X}, Y) is unknown, we compute p-values $\pi_\theta(\mathbf{X}, \mathcal{D})$ depending on the current feature vector \mathbf{X} as well as on the training data \mathcal{D} . In this case, condition (1.7) can be extended in two ways:

$$\mathbb{P}(\pi_\theta(\mathbf{X}, \mathcal{D}) \leq \alpha \mid Y = \theta) \leq \alpha, \quad (1.8)$$

$$\mathbb{P}(\pi_\theta(\mathbf{X}, \mathcal{D}) \leq \alpha \mid Y = \theta, \mathcal{D}) \leq \alpha + o_p(1) \quad \text{as } n \rightarrow \infty, \quad (1.9)$$

for any $\theta \in \mathcal{Y}$ and $\alpha \in (0, 1)$. Condition (1.8) corresponds to “single use” and Condition (1.9) to “multiple use”. Suppose that we construct p-values $\pi_\theta(\cdot, \mathcal{D})$ based on one training data set \mathcal{D} and classify many future observations $(\tilde{\mathbf{X}}, \tilde{Y})$. Then the relative number of future observations with $\tilde{Y} = b$ and $\pi_\theta(\tilde{\mathbf{X}}, \mathcal{D}) \leq \alpha$ is close to

$$w_b \mathbb{P}(\pi_\theta(\mathbf{X}, \mathcal{D}) \leq \alpha \mid Y = b, \mathcal{D}),$$

a random quantity depending on the training data \mathcal{D} .

Example 1.1. For the following one-dimensional example let $L = 2$, $P_1 = \text{Gamma}(3, 1)$ and $P_2 = \text{Gamma}(6, 1)$. Figure 1.1 illustrates how the optimal point predictor $Y^*(x)$ and the posterior weights $w_\theta(x)$ depend on the prior probabilities w_θ . It shows $w_2(x)$ for $w_2/w_1 = 10, 1.5, 1, 0.67, 0.1$ (from left to right). The corresponding boundaries of $Y^*(x)$ are drawn as vertical lines.

Alternatively, we could calculate p-values which do not depend on the prior probabilities w_θ . Since $w_2(x)$ is increasing in x , we define the p-values

$$\pi_1(x) := \mathbb{P}(X \geq x \mid Y = 1),$$

$$\pi_2(x) := \mathbb{P}(X \leq x \mid Y = 2).$$

If $\pi_\theta(X) \leq \alpha$ we claim with confidence $1 - \alpha$ that $Y \neq \theta$. Figure 1.2 shows the p-value functions $\pi_1(x)$ and $\pi_2(x)$. In addition, the three regions where $\hat{\mathcal{Y}}_{0.1}(x) = \{1\}, \{2\}, \{1, 2\}$ are marked.

1.3. Optimal P-Values as Benchmark

In this section we suppose that the distributions P_1, P_2, \dots, P_L have known densities $f_1, f_2, \dots, f_L > 0$ with respect to some measure M on \mathcal{X} . Then the marginal distribution of \mathbf{X} has density

$$f(\mathbf{x}) := \sum_{\theta \in \mathcal{Y}} w_\theta f_\theta(\mathbf{x})$$

1. Classifiers and P-Values

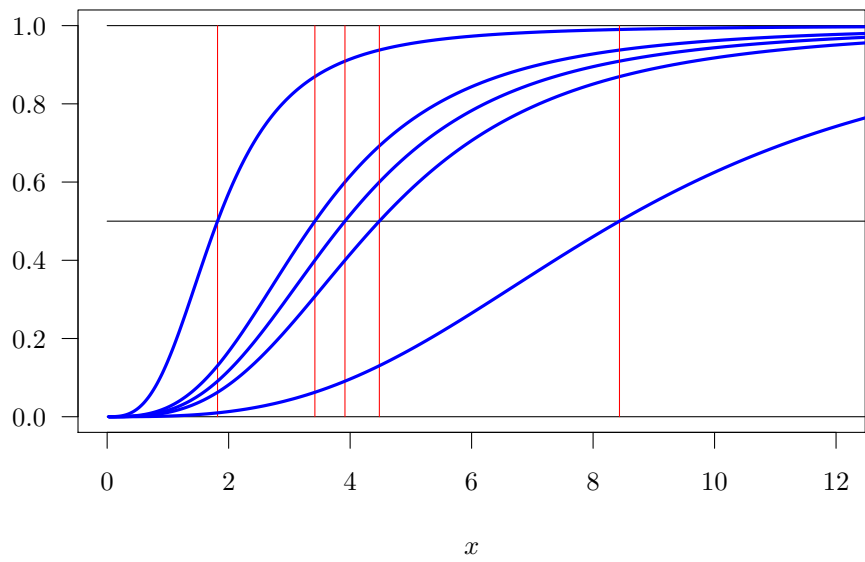


Figure 1.1.: Posterior weights $w_2(x)$ for different ratios of prior probabilities w_2/w_1 .

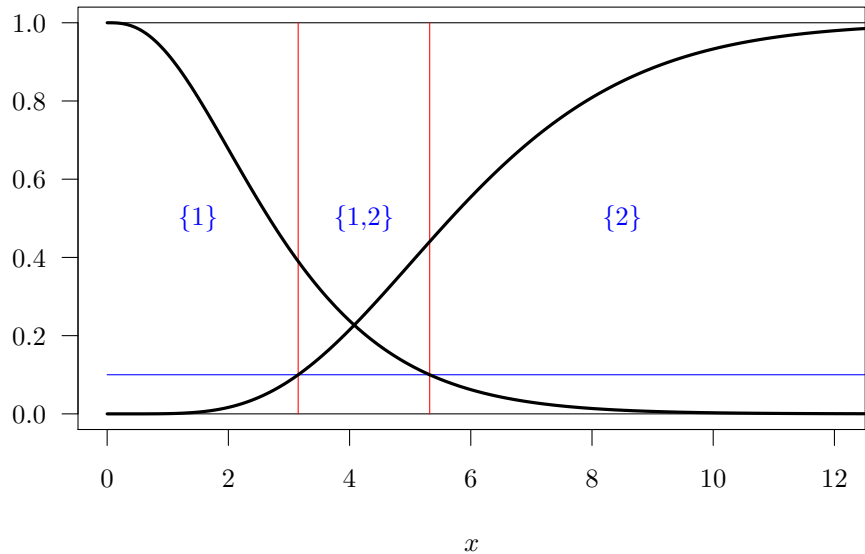


Figure 1.2.: P-value functions for class memberships.

with respect to M and

$$w_\theta(\mathbf{x}) = \frac{w_\theta f_\theta(\mathbf{x})}{f(\mathbf{x})}.$$

Let $\boldsymbol{\pi} = (\pi_\theta)_{\theta \in \mathcal{Y}}$ consist of p-values π_θ satisfying (1.8). Given the latter constraint we want to provide small p-values and small prediction regions. Therefore we use the following measures of risk:

$$\begin{aligned}\mathcal{R}(\boldsymbol{\pi}) &:= \mathbb{E}\left(\sum_{\theta \in \mathcal{Y}} \pi_\theta(\mathbf{X})\right), \\ \mathcal{R}_\alpha(\boldsymbol{\pi}) &:= \mathbb{E}(\#\hat{\mathcal{Y}}_\alpha(\mathbf{X})).\end{aligned}$$

Lemma 1.3. *Define $\mathcal{R}_\alpha(\pi_\theta) := \mathbb{P}(\pi_\theta(\mathbf{X}) > \alpha)$. Then*

$$\mathcal{R}(\boldsymbol{\pi}) = \int_0^1 \mathcal{R}_\alpha(\boldsymbol{\pi}) \, d\alpha$$

and

$$\mathcal{R}_\alpha(\boldsymbol{\pi}) = \sum_{\theta \in \mathcal{Y}} \mathcal{R}_\alpha(\pi_\theta).$$

PROOF. Note that

$$\sum_{\theta \in \mathcal{Y}} \mathcal{R}_\alpha(\pi_\theta) = \mathbb{E}(\#\{\theta \in \mathcal{Y} : \pi_\theta(\mathbf{X}) > \alpha\}) = \mathcal{R}_\alpha(\boldsymbol{\pi})$$

and

$$\mathcal{R}(\boldsymbol{\pi}) = \sum_{\theta \in \mathcal{Y}} \int_0^1 \mathbb{P}(\pi_\theta(\mathbf{X}) > \alpha) \, d\alpha = \int_0^1 \mathcal{R}_\alpha(\boldsymbol{\pi}) \, d\alpha.$$

□

In view of the preceding lemma, we focus on minimizing $\mathcal{R}_\alpha(\pi_\theta)$ for arbitrary fixed $\theta \in \mathcal{Y}$ and $\alpha \in (0, 1)$ under the constraint (1.8).

Lemma 1.4. *Let $\mathcal{L}((f_\theta/f)(\mathbf{X}))$ be continuous. Then the p-value*

$$\pi_\theta^*(\mathbf{x}) := P_\theta\{\mathbf{z} \in \mathcal{X} : (f_\theta/f)(\mathbf{z}) \leq (f_\theta/f)(\mathbf{x})\}$$

is optimal, in the sense that $\mathcal{R}_\alpha(\pi_\theta^)$ is minimal for each $\alpha \in (0, 1)$.*

PROOF. We consider

$$\varphi(\mathbf{x}) := \mathbb{1}\{\pi_\theta(\mathbf{x}) \leq \alpha\}$$

1. Classifiers and P-Values

as a level- α test of the null-hypothesis P_θ versus the alternative hypothesis $P = \sum_{b \in \mathcal{Y}} w_b P_b$ and maximize the power

$$\mathbb{E}\varphi(\mathbf{X}) = \int \varphi(\mathbf{x}) P(d\mathbf{x})$$

subject to the condition

$$\mathbb{E}_\theta \varphi(\mathbf{X}) := \int \varphi(\mathbf{x}) P_\theta(d\mathbf{x}) \leq \alpha.$$

The Neyman-Pearson Lemma (Theorem A.2) yields that $\mathbb{E}\varphi(\mathbf{X})$ is maximal for

$$\varphi(\mathbf{x}) = \begin{cases} 1 & \text{if } (f/f_\theta)(\mathbf{x}) > c_\theta, \\ \gamma_\theta & \text{if } (f/f_\theta)(\mathbf{x}) = c_\theta, \\ 0 & \text{if } (f/f_\theta)(\mathbf{x}) < c_\theta \end{cases}$$

with $c_\theta \in [0, \infty]$ and $\gamma_\theta \in [0, 1]$ such that

$$\begin{aligned} \mathbb{E}_\theta \varphi(\mathbf{X}) &= P_\theta\{\mathbf{x} \in \mathcal{X} : (f/f_\theta)(\mathbf{x}) > c_\theta\} + \gamma_\theta P_\theta\{\mathbf{x} \in \mathcal{X} : (f/f_\theta)(\mathbf{x}) = c_\theta\} \\ &= \alpha. \end{aligned}$$

Since $\mathcal{L}((f_\theta/f)(\mathbf{X}))$ is continuous, γ_θ can be chosen arbitrarily. With $\gamma_\theta = 1$ and

$$c_\theta := \min \left\{ c : P_\theta\{\mathbf{x} \in \mathcal{X} : (f/f_\theta)(\mathbf{x}) \geq c\} \leq \alpha \right\}$$

we get

$$\begin{aligned} \varphi(\mathbf{x}) &= \mathbb{1}\{(f/f_\theta)(\mathbf{x}) \geq c_\theta\} \\ &= \mathbb{1}\{\pi_\theta^*(\mathbf{x}) \leq \alpha\}. \end{aligned}$$

Thus $\varphi(\mathbf{x}) = \mathbb{1}\{\pi_\theta^*(\mathbf{x}) \leq \alpha\}$ maximizes $\mathbb{E}\varphi(\mathbf{x})$ and minimizes $\mathcal{R}_\alpha(\pi_\theta) = 1 - \mathbb{E}\varphi(\mathbf{x})$. \square

Two other representations of $\pi_\theta^*(\mathbf{x})$ are given by

$$\begin{aligned} \pi_\theta^*(\mathbf{x}) &= P_\theta\{\mathbf{z} \in \mathcal{X} : w_\theta(\mathbf{z}) \leq w_\theta(\mathbf{x})\} \\ &= P_\theta\{\mathbf{z} \in \mathcal{X} : T_\theta^*(\mathbf{z}) \geq T_\theta^*(\mathbf{x})\} \end{aligned}$$

with

$$T_\theta^*(\mathbf{x}) := \sum_{b \neq \theta} \frac{w_{b,\theta} f_b(\mathbf{x})}{f_\theta(\mathbf{x})} \quad \text{and} \quad w_{b,\theta} := \left(\sum_{c \neq \theta} w_c / w_b \right)^{-1}.$$

Note that

$$\begin{aligned} T_\theta^*(\cdot) &= \frac{f(\cdot) - w_\theta f_\theta(\cdot)}{f_\theta(\cdot)(1 - w_\theta)} \\ &= \frac{1 - w_\theta(\cdot)}{\frac{f_\theta}{f}(\cdot)(1 - w_\theta)} \\ &= \frac{w_\theta(\cdot)^{-1} - 1}{w_\theta^{-1} - 1} \end{aligned}$$

is a negative monotonic transformation of $w_\theta(\cdot)$. The first representation shows that $\pi_\theta^*(\mathbf{x})$ is a non-decreasing function of $w_\theta(\mathbf{x})$. The second representation shows that the prior weight w_θ itself is irrelevant for the optimal p-value π_θ^* . Only the ratios w_c/w_b with $b, c \neq \theta$ matter. In particular, in case of $L = 2$ classes, the optimal p-values do not depend on the prior distribution of Y at all.

Example 1.2 (Standard model). Let $P_\theta = \mathcal{N}_d(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma})$ with a common covariance matrix $\boldsymbol{\Sigma}$. Then

$$T_\theta^*(\mathbf{x}) := \sum_{b \neq \theta} w_{b,\theta} \exp((\mathbf{x} - \boldsymbol{\mu}_{\theta,b})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_b - \boldsymbol{\mu}_\theta)) \quad (1.10)$$

with $\boldsymbol{\mu}_{\theta,b} = (\boldsymbol{\mu}_\theta + \boldsymbol{\mu}_b)/2$.

1.4. P-Values via Permutation Tests

Now we suppose that the joint distribution of (\mathbf{X}, Y) is unknown and compute p-values $\pi_\theta(\mathbf{X}, \mathcal{D})$ and prediction regions

$$\hat{\mathcal{Y}}_\alpha(\mathbf{X}, \mathcal{D}) := \{\theta \in \mathcal{Y} : \pi_\theta(\mathbf{X}, \mathcal{D}) > \alpha\}$$

depending on training data \mathcal{D} . We introduce nonparametric p-values for classification as they are given in Dümbgen et al. (2008).

For the remaining part of this thesis we consider the class labels Y_1, Y_2, \dots, Y_n as fixed while $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ and (\mathbf{X}, Y) are independent with $\mathcal{L}(\mathbf{X}_i) = P_{Y_i}$. That way we can handle situations with stratified training data as well as the case of i.i.d. training data (via conditioning).

Further we assume that the distributions P_1, P_2, \dots, P_L have densities $f_1, f_2, \dots, f_L > 0$ with respect to some measure M on \mathcal{X} and that all group sizes N_θ are strictly positive. Asymptotic statements as in (1.9) are meant as

$$n \rightarrow \infty \quad \text{and} \quad N_b/n \rightarrow w_b \quad \text{for all } b \in \mathcal{Y}. \quad (1.11)$$

Let $I(1) < I(2) < \dots < I(N_\theta)$ be the elements of \mathcal{G}_θ for a fixed class $\theta \in \mathcal{Y}$. Then $(\mathbf{X}, \mathbf{X}_{I(1)}, \mathbf{X}_{I(2)}, \dots, \mathbf{X}_{I(N_\theta)})$ is exchangeable conditional on $Y = \theta$.

1. Classifiers and P-Values

Thus we consider a test statistic $T_\theta(\mathbf{X}, \mathcal{D})$ which is symmetric in $(\mathbf{X}_{I(j)})_{j=1}^{N_\theta}$. We define $\mathcal{D}_i(\mathbf{x})$ to be the training data with \mathbf{x} in place of \mathbf{X}_i . Then the nonparametric p-value

$$\pi_\theta(\mathbf{X}, \mathcal{D}) := \frac{\#\{i \in \mathcal{G}_\theta : T_\theta(\mathbf{X}_i, \mathcal{D}_i(\mathbf{X})) \geq T_\theta(\mathbf{X}, \mathcal{D})\} + 1}{N_\theta + 1} \quad (1.12)$$

satisfies (1.8). By definition, $\pi_\theta(\mathbf{X}, \mathcal{D}) \geq (N_\theta + 1)^{-1}$. Therefore this procedure is only useful, if $N_\theta + 1 \geq \alpha^{-1}$. For instance if $\alpha = 0.05$, N_θ should be at least 19.

Plug-In Statistic for Standard Model

In the standard model with $P_\theta = \mathcal{N}_d(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma})$ we replace the unknown mean vectors $\boldsymbol{\mu}_\theta$ and covariance matrix $\boldsymbol{\Sigma}$ in (1.10) with corresponding estimators. Note that the resulting p-values always satisfy (1.8), even if the underlying distributions P_θ are not Gaussian with common covariance matrix.

To compute $\pi_\theta(\mathbf{X}, \mathcal{D})$, we add the new observation \mathbf{X} temporarily to class θ . But \mathbf{X} may be an outlier with respect to the distribution P_θ . Therefore it is reasonable to use one of the robust M -estimators mentioned in Section 1.1.2.

Nearest Neighbors

Now we use the test statistic $T_\theta(\mathbf{x}, \mathcal{D}) = -w_\theta(\mathbf{x})$ and estimate $w_\theta(\mathbf{x})$ via nearest neighbors (cf. Section 1.1.2). For the k -nearest neighbors we get

$$\hat{w}_\theta(\mathbf{x}, \mathcal{D}) := \frac{\hat{w}_\theta \hat{P}_\theta(B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x})))}{\sum_{b \in \mathcal{Y}} \hat{w}_b \hat{P}_b(B(\mathbf{x}, \hat{r}_{k,n}(\mathbf{x})))}$$

with certain estimators $\hat{w}_b = \hat{w}_b(\mathcal{D})$ of w_b . In case of $\hat{w}_b = N_b/n$, we get

$$\begin{aligned} \hat{w}_\theta(\mathbf{x}, \mathcal{D}) &:= \frac{\#\{i \in \mathcal{G}_\theta : d(\mathbf{x}, \mathbf{X}_i) \leq \hat{r}_{k,n}(\mathbf{x})\}}{\#\{i \leq n : d(\mathbf{x}, \mathbf{X}_i) \leq \hat{r}_{k,n}(\mathbf{x})\}} \\ &= \frac{\sum_{i=1}^n \mathbb{1}\{d(\mathbf{x}, \mathbf{X}_i) \leq \hat{r}_{k,n}(\mathbf{x})\} \mathbb{1}\{Y_i = \theta\}}{\sum_{i=1}^n \mathbb{1}\{d(\mathbf{x}, \mathbf{X}_i) \leq \hat{r}_{k,n}(\mathbf{x})\}}. \end{aligned}$$

For the weighted nearest neighbors we get

$$\hat{w}_\theta(\mathbf{x}, \mathcal{D}) := \frac{\sum_{i=1}^n W_n(\hat{R}(\mathbf{x}, \mathbf{X}_i)) \mathbb{1}\{Y_i = \theta\}}{\sum_{i=1}^n W_n(\hat{R}(\mathbf{x}, \mathbf{X}_i))}.$$

Penalized Multicategory Logistic Regression

Let $\mathcal{X} = \mathbb{R}^d$ and \mathbf{X} contain the values of d numerical or binary variables. We assume that

$$\mathbb{P}(Y = \theta \mid \mathbf{X} = \mathbf{x}) = \exp(a_\theta + \mathbf{b}_\theta^\top \mathbf{x}) / \sum_{z=1}^L \exp(a_z + \mathbf{b}_z^\top \mathbf{x})$$

for unknown parameters $a_z \in \mathbb{R}$ and $\mathbf{b}_z \in \mathbb{R}^d$, which we estimate with penalized maximum likelihood estimators $\hat{a}_z(\mathcal{D})$ and $\hat{\mathbf{b}}_z(\mathcal{D})$. To compute the p-values, we use the test statistic

$$T_\theta(\mathbf{x}, \mathcal{D}) = -\exp(\hat{a}_\theta + \hat{\mathbf{b}}_\theta^\top \mathbf{x}) / \sum_{z=1}^L \exp(\hat{a}_z + \hat{\mathbf{b}}_z^\top \mathbf{x}).$$

Technical details for the penalized multicategory logistic regression are given in Section 1.8.

1.5. Estimation of Separability

To estimate the separability of different classes by means of given p-values $\pi_\theta(\cdot, \cdot)$ we compute *cross-validated* p-values

$$\pi_\theta(\mathbf{X}_i, \mathcal{D}_i)$$

for $i = 1, 2, \dots, n$ with \mathcal{D}_i denoting the training data without observation (\mathbf{X}_i, Y_i) . We treat each training observation (\mathbf{X}_i, Y_i) temporarily as a 'future' observation, which has to be classified with the remaining data \mathcal{D}_i . Then we could display these p-values graphically or compute the empirical conditional inclusion probabilities

$$\hat{\mathcal{I}}_\alpha(b, \theta) := \frac{\#\{i \in \mathcal{G}_b : \theta \in \hat{\mathcal{Y}}_\alpha(\mathbf{X}_i, \mathcal{D}_i)\}}{N_b}$$

and the empirical pattern probabilities

$$\hat{\mathcal{P}}_\alpha(b, S) := \frac{\#\{i \in \mathcal{G}_b : \hat{\mathcal{Y}}_\alpha(\mathbf{X}_i, \mathcal{D}_i) = S\}}{N_b}$$

for $b, \theta \in \mathcal{Y}$ and $S \subset \mathcal{Y}$. These numbers $\hat{\mathcal{I}}_\alpha(b, \theta)$ and $\hat{\mathcal{P}}_\alpha(b, S)$ can be interpreted as estimators of the conditional inclusion probabilities

$$\mathcal{I}_\alpha(b, \theta \mid \mathcal{D}) := \mathbb{P}(\theta \in \hat{\mathcal{Y}}_\alpha(\mathbf{X}, \mathcal{D}) \mid Y = b, \mathcal{D})$$

1. Classifiers and P-Values

and the conditional pattern probabilities

$$\mathcal{P}_\alpha(b, S \mid \mathcal{D}) := \mathbb{P}(\hat{\mathcal{Y}}_\alpha(\mathbf{X}, \mathcal{D}) = S \mid Y = b, \mathcal{D}),$$

respectively.

To visualize the separability we plot for large group sizes N_b the empirical ROC curves

$$(0, 1) \ni \alpha \mapsto 1 - \hat{\mathcal{I}}_\alpha(b, \theta).$$

1.6. Asymptotic Properties

In this section we review the asymptotic results of Dümbgen et al. (2008) and Zumbrennen (2009). For the plug-in statistic of the standard model we derive a central limit theorem in Section 3.4.

Throughout this section, asymptotic statements are to be understood with-in setting (1.11).

The following theorem implies that $\pi_\theta(\mathbf{X}, \mathcal{D})$ satisfies (1.9) under certain conditions on the underlying test statistic $T_\theta(\mathbf{X}, \mathcal{D})$. Furthermore it shows that the empirical conditional inclusion probabilities $\hat{\mathcal{I}}_\alpha(b, \theta)$ and the empirical pattern probabilities $\hat{\mathcal{P}}_\alpha(b, S)$ are consistent estimators of $\mathcal{I}_\alpha(b, \theta \mid \mathcal{D})$ and $\mathcal{P}_\alpha(b, S \mid \mathcal{D})$, respectively. Here and throughout \rightarrow_p denotes convergence in probability.

Theorem 1.5. *Suppose that for fixed $\theta \in \mathcal{Y}$ there exists a test statistic T_θ^o on \mathcal{X} satisfying the following two requirements:*

$$T_\theta(\mathbf{X}, \mathcal{D}) \rightarrow_p T_\theta^o(\mathbf{X}), \quad (1.13)$$

$$\mathcal{L}(T_\theta^o(\mathbf{X})) \text{ is continuous.} \quad (1.14)$$

Then

$$\pi_\theta(\mathbf{X}, \mathcal{D}) \rightarrow_p \pi_\theta^o(\mathbf{X}), \quad (1.15)$$

where

$$\pi_\theta^o(\mathbf{x}) := P_\theta\{\mathbf{z} \in \mathcal{X} : T_\theta^o(\mathbf{z}) \geq T_\theta^o(\mathbf{x})\}.$$

In particular, for arbitrary fixed $\alpha \in (0, 1)$,

$$\mathcal{R}_\alpha(\pi_\theta(\cdot, \mathcal{D})) \rightarrow_p \mathcal{R}_\alpha(\pi_\theta^o), \quad (1.16)$$

$$\left. \begin{array}{l} \mathcal{I}_\alpha(b, \theta \mid \mathcal{D}) \\ \hat{\mathcal{I}}_\alpha(b, \theta) \end{array} \right\} \rightarrow_p \mathbb{P}(\pi_\theta^o(\mathbf{X}) > \alpha \mid Y = b) \quad \text{for each } b \in \mathcal{Y} \quad (1.17)$$

and

$$\left. \begin{array}{l} \mathcal{P}_\alpha(b, S \mid \mathcal{D}) \\ \hat{\mathcal{P}}_\alpha(b, S) \end{array} \right\} \rightarrow_{\mathbb{P}} \mathbb{P}(\hat{\mathcal{Y}}_\alpha^o(\mathbf{X}) = S) \quad \text{for each } b \in \mathcal{Y} \text{ and } S \subset \mathcal{Y}, \quad (1.18)$$

where $\hat{\mathcal{Y}}_\alpha^o(\mathbf{X}) := \{\theta \in \mathcal{Y} : \pi_\theta^o(\mathbf{X}) > \alpha\}$.

The proof of this theorem can be found in Dümbgen et al. (2008), Theorem 3.1. The p-value $\pi_\theta(\cdot, \mathcal{D})$ is asymptotically optimal if T_θ^o is a strictly increasing transformation of T_θ^* . The following lemmata show that this is the case in different situations.

Plug-In Statistic for Standard Gaussian Model

Lemma 1.6. *Conditions (1.13) and (1.14) are satisfied with $T_\theta^o = T_\theta^*$ in case of the plug-in rule for the homoscedastic Gaussian model, provided that $\mathbb{E}(\|\mathbf{X}\|^2) < \infty$ and $\mathcal{L}(\mathbf{X})$ has a Lebesgue density.*

For the proof we refer to Dümbgen et al. (2008), Lemma 3.2.

Nearest Neighbors

Lemma 1.7. *Suppose that (\mathcal{X}, d) is a separable metric space with a measure M satisfying (1.6) and that all densities f_b , $b \in \mathcal{Y}$, are continuous on \mathcal{X} . Then for the weighted nearest-neighbor rule with weights satisfying*

$$\sum_{i: i \geq \varepsilon n} W_n(i) \rightarrow 0 \quad \text{for all } \varepsilon > 0, \quad (1.19)$$

$$\sum_{i=1}^n W_n(i) = 1 \quad \text{for all } n \in \mathbb{N}, \quad (1.20)$$

$$W_n(1) \rightarrow 0, \quad (1.21)$$

the assumptions of Theorem 1.5 are satisfied with $T_\theta^o = T_\theta^$.*

The proof of this lemma can be found in Zumbunnen (2009), Theorem 3.1. Note that the k -nearest neighbor rule with $\hat{w}_\theta = N_\theta/n$ satisfies the conditions of the previous theorem, provided that

$$k = k(n) \rightarrow \infty \quad \text{and} \quad k/n \rightarrow 0.$$

Often different variables of a data set are measured on different scales. To take this into account, one could use the Mahalanobis distance, which is scale-invariant and data-driven.

1. Classifiers and P-Values

Lemma 1.8. *Let \mathcal{X} be an open subset of \mathbb{R}^d and $f_b, b \in \mathcal{Y}$ continuous Lebesgue densities. Suppose that $\mathbb{E}(\|\mathbf{X}\|^2) < \infty$ and let $\hat{\Sigma}$ be a consistent estimator of the nonsingular matrix $\Sigma_0 := \sum_{\theta \in \Theta} w_\theta \text{Var}(\mathbf{X} \mid Y = \theta)$. Then in case of the weighted nearest-neighbor rule with the Mahalanobis distance $D_{\hat{\Sigma}}$ and weights satisfying (1.19)–(1.21), the assumptions of Theorem 1.5 are satisfied with $T_\theta^o = T_\theta^*$.*

The proof of this lemma is given in Zumbunnen (2009), Theorem 3.10.

1.7. Implementation in `pvclass`

The p-values for classification are implemented in the package `pvclass` (Zumbunnen and Dümngen, 2011). It was written in the R programming system (R Core Team, 2014) and depends on the recommended package `Matrix` (Bates and Maechler, 2010).

The main functions of `pvclass` compute p-values for the potential class memberships of new observations (`pvs`) as well as cross-validated p-values for training data (`cvpvs`). With the function `analyze.pvs`, the package `pvclass` also provides graphical displays and quantitative analyses of the p-values.

The test statistics of Section 1.4 are available in the package `pvclass`. It should be stressed however that users could easily implement test statistics corresponding to their own favorite classifier (e.g. neuronal nets).

To estimate the parameters we use N_θ/n for w_θ and the standard estimator for μ_θ . For Σ the package `pvclass` offers the standard estimator as well as the more robust M -estimators $\hat{\Sigma}_M$ and $\hat{\Sigma}_{sym}$. The estimator $\hat{\Sigma}_{sym}$ requires that the observations \mathbf{X}_i are pairwise different within groups. Otherwise, if an observation occurs more than once, `pvclass` uses only the first to calculate $\hat{\Sigma}_{sym}$.

For the nearest neighbor methods, `pvclass` offers besides the fixed Euclidean distance also two data-driven distances which are scale invariant. The Mahalanobis distance with respect to the estimated covariance matrix $\hat{\Sigma}$ as defined in Section 1.1.1 and the data driven Euclidean distance where we divide each component of \mathbf{X} by its sample standard deviation and then use the Euclidean distance.

For the weighted nearest neighbors `pvclass` provides the linear weight function

$$W_n(i) = \max(1 - (i/n)/\tau, 0),$$

and the exponential weight function

$$W_n(i) = (1 - i/n)^\tau.$$

Alternatively one can specify the weights with an n dimensional vector W . For the exponential weight function τ should be in $(0, 1]$ and for the linear weight function it should be greater than 1.

Details for the test statistic based on penalized multicategory logistic regression are given in Section 1.8.

1.7.1. Shortcut

To reduce computation time, we add (\mathbf{X}, θ) to the training data before we judge the plausibility of the class label θ for a new observation \mathbf{X} , i.e. we replace $T_\theta(\mathbf{X}, \mathcal{D})$ and $T_\theta(\mathbf{X}_i, \mathcal{D}_i(\mathbf{X}))$ in Definition (1.12) by $T_\theta(\mathbf{X}, \mathcal{D}(\mathbf{X}, \theta))$ and $T_\theta(\mathbf{X}_i, \mathcal{D}(\mathbf{X}, \theta))$, respectively, where $\mathcal{D}(\mathbf{X}, \theta)$ denotes the training data extended by (\mathbf{X}, θ) . Then we end up with the p-value

$$\pi_\theta(\mathbf{X}, \mathcal{D}(\mathbf{X}, \theta)) = \frac{\#\{i \in \mathcal{G}_\theta : T_\theta(\mathbf{X}_i, \mathcal{D}(\mathbf{X}, \theta)) \geq T_\theta(\mathbf{X}, \mathcal{D}(\mathbf{X}, \theta))\} + 1}{N_\theta + 1}.$$

To compute $\pi_\theta(\mathbf{X}, \mathcal{D}(\mathbf{X}, \theta))$ it suffices to evaluate $T_\theta(\cdot, \mathcal{D}(\mathbf{X}, \theta))$ at the $N_\theta + 1$ points \mathbf{X} and $\mathbf{X}_i, i \in \mathcal{G}_\theta$. One can show that this p-value satisfies (1.8) and the conclusions of Theorems 1.5 and 3.11 remain true.

1.7.2. Data Example ‘buerk’

To illustrate the main functions of `pvclass` we use the data set `buerk` provided by `pvclass`. It was collected by Prof. Dr. Conny Georg Bürk at the university hospital in Lübeck and contains data of 21’556 surgeries in a certain time period (end of the nineties). Besides the mortality and the morbidity it contains 21 variables describing the condition of the patient and the surgery. All collected variables can be found in Table 1.1.

We use the mortality as class label Y . The original data set contains 21’556 observations. To get a smaller data set, which is easier to handle, we take all 662 observations with $Y = 1$ and choose randomly $3 \cdot 662$ observations with $Y = 0$. For the test data set we choose 100 observations from each class. So we end up with a training data set containing 2448 observations, whereof 562 belong to class 1.

```
R> library(pvclass)
R> data(buerk)
R> set.seed(0)
R> X.raw <- as.matrix(buerk[, 1:21])
R> Y.raw <- buerk[, 22]
R> n0.raw <- sum(1 - Y.raw)
R> n1 <- sum(Y.raw)
R> n0 <- 3 * n1
```

1. Classifiers and P-Values

Variable	Meaning
Y	Mortality (1 = deceased, 0 = survived)
X(1)	Age in years
X(2)	Sex (1 = female, 0 = male)
X(3)	ASA-Score (American Society of Anesthesiologists), describes the physical condition on an ordinal scale 1 = A normal healthy patient, 2 = A patient with mild systemic disease, 3 = A patient with severe systemic disease, 4 = A patient with severe systemic disease that is a constant threat to life, 5 = A moribund patient who is not expected to survive without the operation, 6 = A declared brain-dead patient whose organs are being removed for donor purposes
X(4)	Risk factor: cerebral (1 = yes, 0 = no)
X(5)	Risk factor: cardiovascular (1 = yes, 0 = no)
X(6)	Risk factor: pulmonary (1 = yes, 0 = no)
X(7)	Risk factor: renal (1 = yes, 0 = no)
X(8)	Risk factor: hepatic (1 = yes, 0 = no)
X(9)	Risk factor: immunological (1 = yes, 0 = no)
X(10)	Risk factor: metabolic (1 = yes, 0 = no)
X(11)	Risk factor: uncooperative, unreliable (1 = yes, 0 = no)
X(12)	Etiology: malignant (1 = yes, 0 = no)
X(13)	Etiology: vascular (1 = yes, 0 = no)
X(14)	Antibiotics therapy (1 = yes, 0 = no)
X(15)	Surgery indicated (1 = yes, 0 = no)
X(16)	Emergency operation (1 = yes, 0 = no)
X(17)	Surgery time in minutes
X(18)	Septic surgery (1 = yes, 0 = no)
X(19)	Experienced surgeon, i.e. senior physician (1 = yes, 0 = no)
X(20)	Blood transfusion necessary (1 = yes, 0 = no)
X(21)	Intensive care necessary (1 = yes, 0 = no)

Table 1.1.: Variables in **buerk** data set


```

R> X0 <- X.raw[Y.raw == 0, ]
R> X1 <- X.raw[Y.raw == 1, ]
R> tmpi0 <- sample(1:n0.raw, size = 3 * n1, replace = FALSE)
R> tmpi1 <- sample(1:n1, size = n1, replace = FALSE)
R> Xtrain <- rbind(X0[tmpi0[1:(n0 - 100)], ],
                  X1[1:(n1 - 100), ])
R> Ytrain <- c(rep(0, n0 - 100), rep(1, n1 - 100))
R> Xtest <- rbind(X0[tmpi0[(n0 - 99):n0], ],
                  X1[(n1 - 99):n1, ])
R> Ytest <- c(rep(0, 100), rep(1, 100))

```

1.7.3. Main Functions

Classify new observations

The function `pvs` computes nonparametric p-values for the potential class memberships of new observations. It returns a matrix `PV` containing the p-values. Precisely, for each new observation `NewX[i,]` and each class `b` the number `PV[i,b]` is a p-value for the null hypothesis that $Y[i] = b$. With the option `method` or using directly one of the functions `pvs.method` one can choose a test statistic.

For the following example we use the weighted nearest neighbor statistic with an exponential weight function and `tau = 10`.

```

R> PV <- pvs(NewX = Xtest, X = Xtrain, Y = Ytrain,
             method = 'wnn', wtype = 'exponential', tau = 10)
R> head(PV)

```

```

           0           1
[1,] 0.1738209 0.45470693
[2,] 0.6173821 0.06216696
[3,] 0.1213567 0.58081705
[4,] 0.8473768 0.01776199
[5,] 0.4043455 0.15808171
[6,] 0.2517223 0.34280639

```

Next we illustrate the p-values graphically with the function `analyze.pvs` using the first ten observations of each class.

```

R> analyze.pvs(pv=PV[c(1:10,101:110)],, alpha = 0.05)

```

For each p-value a rectangle with an area proportional to the p-value is drawn, see Figure 1.3. The rectangle is blue if the p-value is greater than `alpha` and red otherwise. If we specify the class labels of the test data as in the next example, then the data are sorted by class and the class labels are

1. Classifiers and P-Values

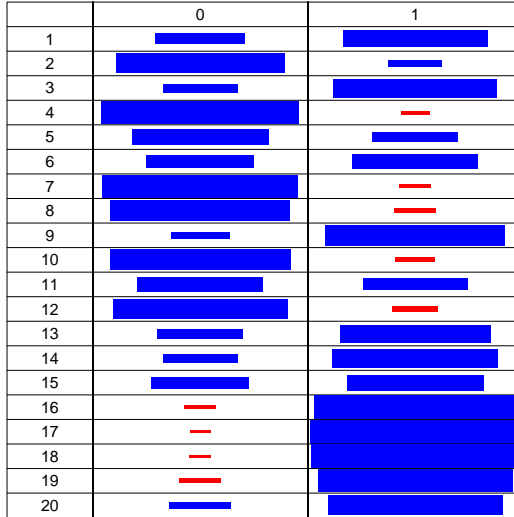


Figure 1.3.: Illustration of the p-values without indicating the class labels of the test data

shown in the plot, see Figure 1.4. Additionally ROC curves are plotted by default. We suppress this here with the argument `roc = FALSE`. An example of the ROC curve plot can be found in the next section.

```
R> analyze.pvs(pv=P[c(1:10,101:110)], ,
+             Y = Ytest[c(1:10,101:110)], roc = FALSE)
```

Cross-validated p-values

The function `cvpvs` returns a matrix `PV` containing cross-validated nonparametric p-values for the potential class memberships of the training data. Precisely, for each feature vector $X[i,]$ and each class b the number $PV[i,b]$ is a p-value for the null hypothesis that $Y[i] = b$.

For the following example we use the logistic regression with penalty parameter `tau.o = 2`.

```
R> PV.cv <- cvpvs(X = Xtrain, Y = Ytrain,
+                 method = 'logreg', tau.o = 2)
R> PV.cv[1:3,]
```

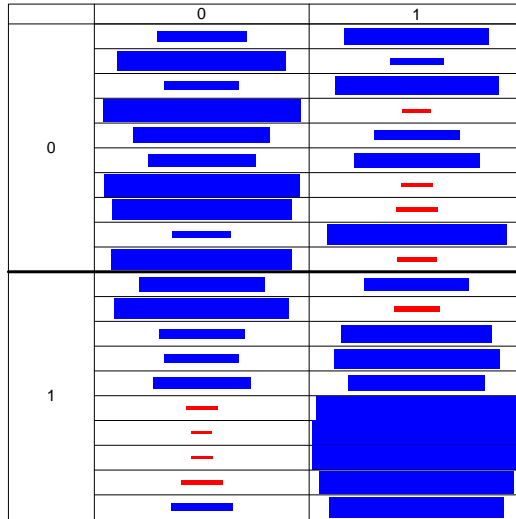


Figure 1.4.: Illustration of the p-values with class labels of the test data

```

      [,1]      [,2]
[1,] 0.9761400 0.001776199
[2,] 0.4172853 0.010657194
[3,] 0.4554613 0.010657194

```

```
R> PV.cv[2001:2003,]
```

```

      [,1]      [,2]
[1,] 0.002119767 0.7971530
[2,] 0.049284579 0.2740214
[3,] 0.010068892 0.6263345

```

The cross-validated p-values can be illustrated graphically the same way as the p-values for the new observations. If $L \leq 3$ the function `analyze.pvs` also prints the empirical pattern probabilities $\hat{P}_\alpha(b, S)$ for all subsets $S \subset \mathcal{Y}$. Otherwise it prints the empirical conditional inclusion probabilities $\hat{I}_\alpha(b, \theta)$ for all combinations of b and θ and the empirical pattern probabilities for $S = \emptyset, \mathcal{Y}$ and $\{\theta\}$ for all $\theta \in \mathcal{Y}$.

In the following example we suppress the plot of the p-values and get only the plot of the ROC curves, see Figure 1.5.

1. Classifiers and P-Values

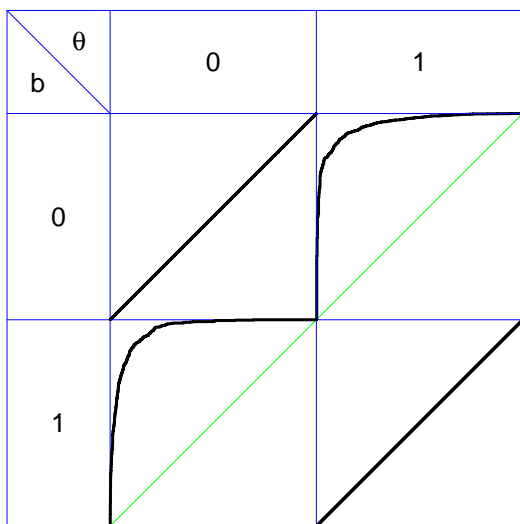


Figure 1.5.: ROC curves of the cross-validated p-values

```
R> analyze.pvs(pv = PV.cv, Y = Ytrain, pvplot = FALSE, cex=1.3)
```

b	P(b,{})	P(b,{1})	P(b,{2})	P(b,{1,2})
1	0	0.78791092	0.04984093	0.1622481
2	0	0.04982206	0.72064057	0.2295374

1.8. Technical Details for Penalized Multicategory Logistic Regression

One of our versions of penalized multicategory logistic regression is similar to the regularized multinomial regression introduced by Friedman et al. (2010), the other one is a variation of the procedure of Zhu and Hastie (2004). An important difference is that we use a smooth approximation to the absolute value or norm function so that usual Newton-Raphson procedures (with step size correction) are applicable.

1.8.1. The Log-Likelihood-Function

Let (\mathbf{X}, Y) be a random variable with values in $\mathbb{R}^d \times \{1, \dots, L\}$ for some integer $L \geq 2$. We assume that

$$P(Y = y | \mathbf{X} = \mathbf{x}) = \exp(a_y + \mathbf{b}_y^\top \mathbf{x}) / \sum_{z=1}^L \exp(a_z + \mathbf{b}_z^\top \mathbf{x})$$

for unknown parameters $a_y \in \mathbb{R}$ and $\mathbf{b}_y \in \mathbb{R}^d$. For notational convenience we introduce the vectors

$$\mathbf{V} = \begin{pmatrix} 1 \\ \mathbf{X} \end{pmatrix} \quad \text{and} \quad \boldsymbol{\theta}_y = \begin{pmatrix} a_y \\ \mathbf{b}_y \end{pmatrix}$$

in $\mathbb{R}^{d'}$ with $d' := 1 + d$. Then

$$P(Y = y | \mathbf{V} = \mathbf{v}) = \exp(\boldsymbol{\theta}_y^\top \mathbf{v}) / \sum_{z=1}^L \exp(\boldsymbol{\theta}_z^\top \mathbf{v}).$$

This parametrization is not unique, because $P(Y = y | \mathbf{V} = \mathbf{v})$ remains unchanged if we add one and the same arbitrary vector to all parameters $\boldsymbol{\theta}_y$. We will deal with this non-uniqueness later in various ways. Our goal is estimation of

$$\boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \boldsymbol{\theta}_2^\top, \dots, \boldsymbol{\theta}_L^\top)^\top \in \mathbb{R}^{Ld'},$$

1. Classifiers and P-Values

based on independent data pairs (\mathbf{V}_i, Y_i) , $1 \leq i \leq n$, such that $\mathcal{L}(Y_i | \mathbf{V}_i) = \mathcal{L}(Y | \mathbf{V})$. Thus we consider the negative log-likelihood function

$$\Lambda(\boldsymbol{\theta}) := \sum_{i=1}^n \left(-\boldsymbol{\theta}_{Y_i}^\top \mathbf{V}_i + \log \left(\sum_{y=1}^L \exp(\boldsymbol{\theta}_y^\top \mathbf{V}_i) \right) \right).$$

For the computation of the first and second derivatives of $\Lambda(\cdot)$ we shall use the following lemma:

Lemma 1.9. *Consider the functional*

$$\mathbb{R}^L \ni \mathbf{f} \mapsto \lambda(\mathbf{f}) := \log \left(\sum_{y=1}^L \exp(f_y) \right).$$

The gradient vector and Hessian matrix of this functional λ at \mathbf{f} are given by

$$\mathbf{p}(\mathbf{f}) := \left(\exp(f_y) / \sum_{z=1}^L \exp(f_z) \right)_{y=1}^L$$

and

$$\mathbf{h}(\mathbf{f}) := \text{diag}(\mathbf{p}(\mathbf{f})) - \mathbf{p}(\mathbf{f})\mathbf{p}(\mathbf{f})^\top,$$

respectively. Moreover, for any $\mathbf{v} \in \mathbb{R}^L$,

$$\mathbf{v}^\top \mathbf{h}(\mathbf{f}) \mathbf{v} \begin{cases} = 0 & \text{if } v_1 = v_2 = \cdots = v_L, \\ > 0 & \text{else.} \end{cases}$$

PROOF. The formulae for gradient vector and Hessian matrix follow from elementary calculations. As to the sign of $\mathbf{v}^\top \mathbf{h}(\mathbf{f}) \mathbf{v}$, note that

$$\mathbf{v}^\top \mathbf{h}(\mathbf{f}) \mathbf{v} = \sum_{y=1}^L p_y(\mathbf{f}) v_y^2 - \left(\sum_{z=1}^L p_z(\mathbf{f}) v_z \right)^2 = \sum_{y=1}^L p_y(\mathbf{f}) (v_y - \bar{v}(\mathbf{f}))^2,$$

where $\bar{v}(\mathbf{f})$ stands for the weighted average $\sum_{y=1}^L p_y(\mathbf{f}) v_y$. Thus $\mathbf{v}^\top \mathbf{h}(\mathbf{f}) \mathbf{v}$ is non-negative and equals zero if, and only if, all components of \mathbf{v} are identical. \square

With Lemma 1.9 at hand one can easily determine the first and second derivatives of $\Lambda(\cdot)$. To formulate the results we use the Kronecker product

$\mathbf{B} \otimes \mathbf{C}$ of arbitrary matrices (or vectors) \mathbf{B} and \mathbf{C} , namely

$$\mathbf{B} \otimes \mathbf{C} := \begin{pmatrix} B_{11}\mathbf{C} & B_{12}\mathbf{C} & B_{13}\mathbf{C} & \cdots \\ B_{21}\mathbf{C} & B_{22}\mathbf{C} & B_{23}\mathbf{C} & \cdots \\ B_{31}\mathbf{C} & B_{32}\mathbf{C} & B_{33}\mathbf{C} & \cdots \\ \vdots & \vdots & \vdots & \end{pmatrix}.$$

For our purposes it is useful to know that

$$(\mathbf{B} \otimes \mathbf{C})^\top = \mathbf{B}^\top \otimes \mathbf{C}^\top \quad (1.22)$$

and

$$(\mathbf{B} \otimes \mathbf{C})(\mathbf{D} \otimes \mathbf{E}) = (\mathbf{B}\mathbf{D}) \otimes (\mathbf{C}\mathbf{E}) \quad (1.23)$$

for arbitrary matrices $\mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{E}$ such that $\mathbf{B}\mathbf{D}$ and $\mathbf{C}\mathbf{E}$ are well-defined.

Moreover, for any dimension q , the standard basis of \mathbb{R}^q is denoted by $\mathbf{e}_{q,1}, \mathbf{e}_{q,2}, \dots, \mathbf{e}_{q,q}$.

Theorem 1.10. *Let $\mathbf{f}(\boldsymbol{\theta}, \mathbf{v}) := (\boldsymbol{\theta}_y^\top \mathbf{v})_{y=1}^L$ for $\mathbf{v} \in \mathbb{R}^{d'}$. With $\mathbf{p}(\cdot)$ and $\mathbf{h}(\cdot)$ as in Lemma 1.9, the gradient vector and Hessian matrix of the negative log-likelihood $\Lambda(\cdot)$ at $\boldsymbol{\theta}$ are given by*

$$\mathbf{G}(\boldsymbol{\theta}) = \sum_{i=1}^n (\mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) - \mathbf{e}_{L, Y_i}) \otimes \mathbf{V}_i$$

and

$$\mathbf{H}(\boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \otimes (\mathbf{V}_i \mathbf{V}_i^\top),$$

respectively. The matrix $\mathbf{H}(\boldsymbol{\theta})$ is positive semidefinite. If the linear span of $\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_n$ equals $\mathbb{R}^{d'}$, then for arbitrary $\boldsymbol{\delta} = (\boldsymbol{\delta}_1^\top, \boldsymbol{\delta}_2^\top, \dots, \boldsymbol{\delta}_L^\top)^\top \in \mathbb{R}^{Ld'}$,

$$\boldsymbol{\delta}^\top \mathbf{H}(\boldsymbol{\theta}) \boldsymbol{\delta} = 0 \quad \text{if, and only if,} \quad \boldsymbol{\delta}_1 = \boldsymbol{\delta}_2 = \dots = \boldsymbol{\delta}_L.$$

PROOF. Since $\mathbf{f}(\boldsymbol{\theta}, \mathbf{v})$ is linear in $\boldsymbol{\theta}$, it follows from Lemma 1.9 that

$$\begin{aligned} & \lambda(\mathbf{f}(\boldsymbol{\theta} + \boldsymbol{\delta}, \mathbf{v})) - \lambda(\mathbf{f}(\boldsymbol{\theta}, \mathbf{v})) \\ &= \lambda(\mathbf{f}(\boldsymbol{\theta}, \mathbf{v}) + \mathbf{f}(\boldsymbol{\delta}, \mathbf{v})) - \lambda(\mathbf{f}(\boldsymbol{\theta}, \mathbf{v})) \\ &= \mathbf{f}(\boldsymbol{\delta}, \mathbf{v})^\top \mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{v})) + 2^{-1} \mathbf{f}(\boldsymbol{\delta}, \mathbf{v})^\top \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{v})) \mathbf{f}(\boldsymbol{\delta}, \mathbf{v}) + o(\|\boldsymbol{\delta}\|^2) \end{aligned}$$

as $\boldsymbol{\delta} \rightarrow \mathbf{0}$. Consequently,

$$\begin{aligned} \Lambda(\boldsymbol{\theta} + \boldsymbol{\delta}) - \Lambda(\boldsymbol{\theta}) &= \sum_{i=1}^n (-\boldsymbol{\delta}_{Y_i}^\top \mathbf{V}_i + \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i)^\top \mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i))) \\ &\quad + 2^{-1} \sum_{i=1}^n \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i)^\top \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i) + o(\|\boldsymbol{\delta}\|^2). \end{aligned}$$

1. Classifiers and P-Values

To obtain gradient and Hessian of $\Lambda(\cdot)$ explicitly, note that $\mathbf{f}(\boldsymbol{\theta}, \mathbf{v}) = (\mathbb{I}_L \otimes \mathbf{v})^\top \boldsymbol{\theta} = (\mathbb{I}_L \otimes \mathbf{v}^\top) \boldsymbol{\theta}$. Thus the linear term (in $\boldsymbol{\delta}$) of the previous expansion of $\Lambda(\boldsymbol{\theta} + \boldsymbol{\delta})$ equals

$$\begin{aligned}
& \sum_{i=1}^n (\mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i)^\top \mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) - \boldsymbol{\delta}_{Y_i}^\top \mathbf{V}_i) \\
&= \boldsymbol{\delta}^\top \sum_{i=1}^n ((\mathbb{I}_L \otimes \mathbf{V}_i) \mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) - \mathbf{e}_{L, Y_i} \otimes \mathbf{V}_i) \\
&= \boldsymbol{\delta}^\top \sum_{i=1}^n ((\mathbb{I}_L \otimes \mathbf{V}_i) (\mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \otimes 1) - \mathbf{e}_{L, Y_i} \otimes \mathbf{V}_i) \\
&= \boldsymbol{\delta}^\top \sum_{i=1}^n (\mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \otimes \mathbf{V}_i - \mathbf{e}_{L, Y_i} \otimes \mathbf{V}_i) \\
&= \boldsymbol{\delta}^\top \sum_{i=1}^n (\mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) - \mathbf{e}_{L, Y_i}) \otimes \mathbf{V}_i,
\end{aligned}$$

while twice the quadratic term may be written as

$$\begin{aligned}
& \sum_{i=1}^n \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i)^\top \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i) \\
&= \boldsymbol{\delta}^\top \sum_{i=1}^n (\mathbb{I}_L \otimes \mathbf{V}_i) \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) (\mathbb{I}_L \otimes \mathbf{V}_i^\top) \boldsymbol{\delta} \\
&= \boldsymbol{\delta}^\top \sum_{i=1}^n (\mathbb{I}_L \otimes \mathbf{V}_i) (\mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \otimes 1) (\mathbb{I}_L \otimes \mathbf{V}_i^\top) \boldsymbol{\delta} \\
&= \boldsymbol{\delta}^\top \left(\sum_{i=1}^n \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \otimes (\mathbf{V}_i \mathbf{V}_i^\top) \right) \boldsymbol{\delta}.
\end{aligned}$$

Finally, note that

$$\boldsymbol{\delta}^\top \mathbf{H}(\boldsymbol{\theta}) \boldsymbol{\delta} = \sum_{i=1}^n \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i)^\top \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i) \geq 0$$

with equality if, and only if, $\mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i)^\top \mathbf{h}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) \mathbf{f}(\boldsymbol{\delta}, \mathbf{V}_i) = 0$ for all indices i . According to Lemma 1.9, the latter condition is equivalent to

$$\boldsymbol{\delta}_1^\top \mathbf{V}_i = \boldsymbol{\delta}_2^\top \mathbf{V}_i = \cdots = \boldsymbol{\delta}_L^\top \mathbf{V}_i \quad \text{for } 1 \leq i \leq n.$$

But if the vectors \mathbf{V}_i span the whole $\mathbb{R}^{d'}$, this is equivalent to all vectors $\boldsymbol{\delta}_y$ being identical. \square

1.8.2. Regularizations

Regularization 0. One way to guarantee uniqueness of the parameter $\boldsymbol{\theta}$ is to require

$$\sum_{y=1}^L \boldsymbol{\theta}_y = \mathbf{0}.$$

More generally, let $\boldsymbol{\theta}_{[j]} := (\theta_{j,y})_{y=1}^L$ with $\theta_{j,y}$ denoting the j -th component of $\boldsymbol{\theta}_y$. With $\mathbf{1}_L := (1, 1, \dots, 1)^\top \in \mathbb{R}^L$, the previous condition means that

$$\mathbf{1}_L^\top \boldsymbol{\theta}_{[j]} = 0 \quad (1.24)$$

for all $j = 1, 2, \dots, d'$. To enforce (1.24) at least for some j we can add

$$R_0(\boldsymbol{\theta}) := 2^{-1} \sum_{j=1}^{d'} \sigma_j (\mathbf{1}_L^\top \boldsymbol{\theta}_{[j]})^2$$

with a vector $\boldsymbol{\sigma} = (\sigma_j)_{j=1}^{d'} \in [0, \infty)^{d'}$ to $\Lambda(\boldsymbol{\theta})$. The choice of $\boldsymbol{\sigma}$ will depend on further regularization terms.

Theorem 1.11. *The gradient vector and Hessian matrix of R_0 at $\boldsymbol{\theta}$ are given by*

$\mathbf{G}_{R,0}(\boldsymbol{\theta}) = ((\mathbf{1}_L \mathbf{1}_L^\top) \otimes \text{diag}(\boldsymbol{\sigma})) \boldsymbol{\theta}$ and $\mathbf{H}_{R,0}(\boldsymbol{\theta}) = (\mathbf{1}_L \mathbf{1}_L^\top) \otimes \text{diag}(\boldsymbol{\sigma})$, respectively.

PROOF. Expanding $R_0(\cdot)$ is rather simple, because it is a quadratic functional itself. Note first that

$$\begin{aligned} R_0(\boldsymbol{\theta} + \boldsymbol{\delta}) - R_0(\boldsymbol{\theta}) &= \sum_{j=1}^{d'} \sigma_j (\mathbf{1}_L^\top \boldsymbol{\delta}_{[j]}) (\mathbf{1}_L^\top \boldsymbol{\theta}_{[j]}) + 2^{-1} \sum_{j=1}^{d'} \sigma_j (\mathbf{1}_L^\top \boldsymbol{\delta}_{[j]})^2 \\ &= \sum_{j=1}^{d'} \sigma_j \boldsymbol{\delta}_{[j]}^\top \mathbf{1}_L \mathbf{1}_L^\top \boldsymbol{\theta}_{[j]} + 2^{-1} \sum_{j=1}^{d'} \sigma_j \boldsymbol{\delta}_{[j]}^\top \mathbf{1}_L \mathbf{1}_L^\top \boldsymbol{\delta}_{[j]}. \end{aligned}$$

But the subvector $\mathbf{a}_{[j]}$ of $\mathbf{a} \in \{\boldsymbol{\theta}, \boldsymbol{\delta}\}$ may be written as $(\mathbb{I}_L \otimes \mathbf{e}_{d',j}^\top) \mathbf{a}$, so

$$\begin{aligned} \boldsymbol{\delta}_{[j]}^\top \mathbf{1}_L \mathbf{1}_L^\top \mathbf{a}_{[j]} &= \boldsymbol{\delta}^\top (\mathbb{I}_L \otimes \mathbf{e}_{d',j}) \mathbf{1}_L \mathbf{1}_L^\top (\mathbb{I}_L \otimes \mathbf{e}_{d',j}^\top) \mathbf{a} \\ &= \boldsymbol{\delta}^\top (\mathbb{I}_L \otimes \mathbf{e}_{d',j}) ((\mathbf{1}_L \mathbf{1}_L^\top) \otimes \mathbf{1}) (\mathbb{I}_L \otimes \mathbf{e}_{d',j}^\top) \mathbf{a} \\ &= \boldsymbol{\delta}^\top ((\mathbf{1}_L \mathbf{1}_L^\top) \otimes (\mathbf{e}_{d',j} \mathbf{e}_{d',j}^\top)) \mathbf{a}, \end{aligned}$$

and $\sum_{j=1}^{d'} \sigma_j (\mathbf{1}_L \mathbf{1}_L^\top) \otimes (\mathbf{e}_{d',j} \mathbf{e}_{d',j}^\top) = (\mathbf{1}_L \mathbf{1}_L^\top) \otimes \text{diag}(\boldsymbol{\sigma})$. Hence

$$R_0(\boldsymbol{\theta} + \boldsymbol{\delta}) - R_0(\boldsymbol{\theta}) = \boldsymbol{\delta}^\top ((\mathbf{1}_L \mathbf{1}_L^\top) \otimes \text{diag}(\boldsymbol{\sigma})) \boldsymbol{\theta} + 2^{-1} \boldsymbol{\delta}^\top ((\mathbf{1}_L \mathbf{1}_L^\top) \otimes \text{diag}(\boldsymbol{\sigma})) \boldsymbol{\delta}.$$

□

Regularization 1: penalizing subvectors. For logistic regression there are various good reasons to regularize the functional Λ or $\Lambda + R_0$. One is to avoid numerical problems. Another is to guarantee existence of a minimizer in cases where Λ alone has no minimizer. This happens if one subgroup $\{\mathbf{X}_i: Y_i = \theta_o\}$ is separated from $\{\mathbf{X}_i: Y_i \neq \theta_o\}$ by a hyperplane. Moreover, we want to favor parameter vectors with only few large components. A first way to do this would be to add the penalty

$$\sum_{j=1}^{d'} \tau_j \|\boldsymbol{\theta}_{[j]}\|$$

with $\boldsymbol{\tau} = (\tau_j)_{j=1}^{d'} \in [0, \infty)^{d'}$ to $\Lambda(\boldsymbol{\theta}) + R_0(\boldsymbol{\theta})$. Here and throughout, $\|\cdot\|$ denotes Euclidean norm. This regularization is motivated by Tibshirani's (1996) LASSO and similar in spirit to penalized logistic regression as proposed by Zhu and Hastie (2004). The latter authors used $\|\boldsymbol{\theta}_{[j]}\|^2$ instead of $\|\boldsymbol{\theta}_{[j]}\|$. To avoid problems with the non-smoothness of $\|\cdot\|$ at zero, we approximate it by a smooth function and consider

$$R_1(\boldsymbol{\theta}) := \sum_{j=1}^{d'} \tau_j (\varepsilon^2 + \|\boldsymbol{\theta}_{[j]}\|^2)^{1/2}$$

for some small number $\varepsilon > 0$. Typically we consider $\tau_1 = 0$ and strictly positive parameters $\tau_2, \dots, \tau_{d'}$. Note that $\|\boldsymbol{\theta}_{[j]} - c \mathbf{1}_L\|^2$ becomes minimal if c equals the mean $\mathbf{1}_L^\top \boldsymbol{\theta}_{[j]} / L$. Hence minimizing $\Lambda(\boldsymbol{\theta}) + R_0(\boldsymbol{\theta}) + R_1(\boldsymbol{\theta})$ enforces Condition (1.24) whenever $\sigma_j + \tau_j > 0$.

The following lemma is useful for the analysis of R_1 :

Lemma 1.12. *Consider the functional*

$$\mathbb{R}^L \ni \mathbf{f} \mapsto \rho(\mathbf{f}) := (\varepsilon^2 + \|\mathbf{f}\|^2)^{1/2}.$$

The gradient vector and Hessian matrix of this functional ρ at \mathbf{f} are given by

$$\mathbf{g}_\rho(\mathbf{f}) := \rho(\mathbf{f})^{-1} \mathbf{f} \quad \text{and} \quad \mathbf{h}_\rho(\mathbf{f}) := \rho(\mathbf{f})^{-1} \mathbb{I}_L - \rho(\mathbf{f})^{-3} \mathbf{f} \mathbf{f}^\top,$$

respectively. Moreover, $\mathbf{h}_\rho(\mathbf{f})$ is positive definite for any $\mathbf{f} \in \mathbb{R}^L$.

PROOF. Since $(1 + \delta)^{1/2} = 1 + \delta/2 - \delta^2/8 + O(\delta^3)$ as $\delta \rightarrow 0$,

$$\begin{aligned}
 \rho(\mathbf{f} + \mathbf{v}) &= (\varepsilon^2 + \|\mathbf{f}\|^2 + 2\mathbf{f}^\top \mathbf{v} + \|\mathbf{v}\|^2)^{1/2} \\
 &= \rho(\mathbf{f}) \left(1 + \rho(\mathbf{f})^{-2} (2\mathbf{f}^\top \mathbf{v} + \|\mathbf{v}\|^2) \right)^{1/2} \\
 &= \rho(\mathbf{f}) + \rho(\mathbf{f})^{-1} (2\mathbf{f}^\top \mathbf{v} + \|\mathbf{v}\|^2)/2 - \rho(\mathbf{f})^{-3} (2\mathbf{f}^\top \mathbf{v} + \|\mathbf{v}\|^2)^2/8 \\
 &\quad + O(\|\mathbf{v}\|^3) \\
 &= \rho(\mathbf{f}) + \rho(\mathbf{f})^{-1} \mathbf{f}^\top \mathbf{v} + \mathbf{v}^\top (\rho(\mathbf{f})^{-1} \mathbb{I} - \rho(\mathbf{f})^{-3} \mathbf{f} \mathbf{f}^\top) \mathbf{v} / 2 + O(\|\mathbf{v}\|^3).
 \end{aligned}$$

This proves that gradient and Hessian of ρ at \mathbf{f} are given by $\mathbf{g}_\rho(\mathbf{f}) := \rho(\mathbf{f})^{-1} \mathbf{f}$ and $\mathbf{h}_\rho(\mathbf{f}) := \rho(\mathbf{f})^{-1} \mathbb{I} - \rho(\mathbf{f})^{-3} \mathbf{f} \mathbf{f}^\top$, respectively. Moreover, since $\rho(\mathbf{f}) > \|\mathbf{f}\|$, it follows from the Cauchy-Schwarz inequality that

$$\mathbf{v}^\top \mathbf{h}_\rho(\mathbf{f}) \mathbf{v} \geq \rho(\mathbf{f})^{-1} (1 - \rho(\mathbf{f})^{-2} \|\mathbf{f}\|^2) \|\mathbf{v}\|^2,$$

which is strictly positive for $\mathbf{v} \neq 0$. \square

By means of Lemma 1.12 one can determine the first and second derivatives of the regularizing function $R_1(\cdot)$:

Theorem 1.13. *With $\mathbf{g}_\rho(\cdot)$ and $\mathbf{h}_\rho(\cdot)$ as in Lemma 1.12, the gradient vector and Hessian matrix of $R_1(\cdot)$ at $\boldsymbol{\theta}$ are given by*

$$\mathbf{G}_{R,1}(\boldsymbol{\theta}) = \sum_{j=1}^{d'} \tau_j (\mathbf{g}_\rho(\boldsymbol{\theta}_{[j]}) \otimes \mathbf{e}_{d',j})$$

and

$$\mathbf{H}_{R,1}(\boldsymbol{\theta}) = \sum_{j=1}^{d'} \tau_j (\mathbf{h}_\rho(\boldsymbol{\theta}_{[j]}) \otimes (\mathbf{e}_{d',j} \mathbf{e}_{d',j}^\top)),$$

respectively. Moreover, for any $\boldsymbol{\delta} = (\boldsymbol{\delta}_1^\top, \dots, \boldsymbol{\delta}_L^\top)^\top$ in $\mathbb{R}^{Ld'}$, $\boldsymbol{\delta}^\top \mathbf{H}_{R,1}(\boldsymbol{\theta}) \boldsymbol{\delta} \geq 0$ with equality if, and only if,

$$\tau_1 \boldsymbol{\delta}_{[1]} = \tau_2 \boldsymbol{\delta}_{[2]} = \dots = \tau_{d'} \boldsymbol{\delta}_{[d']} = \mathbf{0}.$$

PROOF. It follows from Lemma 1.12 that

$$\begin{aligned}
 R_1(\boldsymbol{\theta} + \boldsymbol{\delta}) - R_1(\boldsymbol{\theta}) &= \sum_{j=1}^{d'} \tau_j \boldsymbol{\delta}_{[j]}^\top \mathbf{g}_\rho(\boldsymbol{\theta}_{[j]}) + 2^{-1} \sum_{j=1}^{d'} \tau_j \boldsymbol{\delta}_{[j]}^\top \mathbf{h}_\rho(\boldsymbol{\theta}_{[j]}) \boldsymbol{\delta}_{[j]} + o(\|\boldsymbol{\delta}\|^2).
 \end{aligned} \tag{1.25}$$

1. Classifiers and P-Values

But $\boldsymbol{\delta}_{[j]} = (\mathbb{I}_L \otimes \mathbf{e}_{d',j})^\top \boldsymbol{\delta}$, whence

$$\begin{aligned} \boldsymbol{\delta}_{[j]}^\top \mathbf{g}_\rho(\boldsymbol{\theta}_{[j]}) &= \boldsymbol{\delta}^\top (\mathbb{I}_L \otimes \mathbf{e}_{d',j}) (\mathbf{g}_\rho(\boldsymbol{\theta}_{[j]}) \otimes 1) \\ &= \boldsymbol{\delta}^\top (\mathbf{g}_\rho(\boldsymbol{\theta}_{[j]}) \otimes \mathbf{e}_{d',j}), \\ \boldsymbol{\delta}_{[j]}^\top \mathbf{h}_\rho(\boldsymbol{\theta}_{[j]}) \boldsymbol{\delta}_{[j]} &= \boldsymbol{\delta}^\top (\mathbb{I}_L \otimes \mathbf{e}_{d',j}) (\mathbf{h}_\rho(\boldsymbol{\theta}_{[j]}) \otimes 1) (\mathbb{I}_L \otimes \mathbf{e}_{d',j}^\top) \boldsymbol{\delta} \\ &= \boldsymbol{\delta}^\top (\mathbf{h}_\rho(\boldsymbol{\theta}_{[j]}) \otimes (\mathbf{e}_{d',j} \mathbf{e}_{d',j}^\top)) \boldsymbol{\delta}. \end{aligned}$$

Plugging in the previous expressions within (1.25) yields the asserted expressions for gradient and Hessian.

The additional assertion about the Hessian matrix $\mathbf{H}_{R,1}(\boldsymbol{\theta})$ follows from (1.25) and the fact that all matrices $\mathbf{h}_\rho(\boldsymbol{\theta}_{[j]})$ are positive definite. \square

Regularization 2: component-wise penalties. A simple form of regularization, analogous to Tibshirani's (1996) LASSO is to add the penalty

$$\sum_{j=1}^{d'} \tau_j \sum_{y=1}^L |\boldsymbol{\theta}_{j,y}| = \sum_{j=1}^{d'} \tau_j \|\boldsymbol{\theta}_{[j]}\|_1$$

to $\Lambda(\boldsymbol{\theta}) + R_0(\boldsymbol{\theta})$. Again we use a smoothed version of this, namely

$$R_2(\boldsymbol{\theta}) := \sum_{j=1}^{d'} \tau_j \sum_{y=1}^L (\varepsilon^2 + \theta_{j,y}^2)^{1/2}.$$

Application of Lemma 1.12 in the special case of $L = 1$ yields the derivatives

$$\rho'(f) = \frac{f}{(\varepsilon^2 + f^2)^{1/2}} \quad \text{and} \quad \rho''(f) = \frac{\varepsilon^2}{(\varepsilon^2 + f^2)^{3/2}}$$

of $\mathbb{R} \ni f \mapsto \rho(f) = (\varepsilon^2 + f^2)^{1/2}$. Hence the first two derivatives of R_2 have a rather simple form:

Theorem 1.14. *Let $\mathbf{v} = \mathbf{v}(\boldsymbol{\tau}, \boldsymbol{\theta}) \in \mathbb{R}^{Ld'}$ contain the vectors $\text{diag}(\boldsymbol{\tau})\boldsymbol{\theta}_1, \text{diag}(\boldsymbol{\tau})\boldsymbol{\theta}_2, \dots, \text{diag}(\boldsymbol{\tau})\boldsymbol{\theta}_L$ from top to bottom. The gradient vector and Hessian matrix of $R_2(\cdot)$ at $\boldsymbol{\theta}$ are given by*

$$\mathbf{G}_{R,2}(\boldsymbol{\theta}) = \text{vec}((\tau_j \rho'(\theta_{j,y}))_{j \leq d', y \leq L})$$

and

$$\mathbf{H}_{R,2}(\boldsymbol{\theta}) = \text{diag}\left(\text{vec}((\tau_j \rho''(\theta_{j,y}))_{j \leq d', y \leq L})\right),$$

respectively.

In the previous theorem we use the notation $\text{vec}(\mathbf{M})$ for a vector which is formed by stacking the columns of a matrix \mathbf{M} (from left to right).

1.8.3. Strict Convexity and Coercivity

In this subsection we tackle the question when a unique minimizer of $\Lambda + R_0$ or of $\Lambda + R$ with $R = R_0 + R_1$ or $R_0 + R_2$ exists. Let us start with a general consideration: Suppose that $f : \mathbb{R}^q \rightarrow \mathbb{R}$ is continuously differentiable and convex. Then one can easily verify that the following three statements are equivalent:

$$\text{The set of minimizers of } f \text{ is nonvoid and compact;} \quad (1.26)$$

$$f \text{ is coercive, i.e. } f(\boldsymbol{\theta}) \rightarrow \infty \text{ as } \|\boldsymbol{\theta}\| \rightarrow \infty; \quad (1.27)$$

$$\lim_{t \rightarrow \infty} \boldsymbol{\theta}^\top \nabla f(t\boldsymbol{\theta}) > 0 \text{ for any } \boldsymbol{\theta} \in \mathbb{R}^q \setminus \{\mathbf{0}\}. \quad (1.28)$$

The third statement (1.28) becomes more plausible when noting that $\mathbb{R} \ni t \mapsto \boldsymbol{\theta}^\top \nabla f(t\boldsymbol{\theta})$ is the derivative of the convex function $\mathbb{R} \ni t \mapsto f(t\boldsymbol{\theta})$.

Theorem 1.15. *Suppose that $\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_n$ span $\mathbb{R}^{d'}$, and let $\boldsymbol{\sigma} \in (0, \infty)^{d'}$. Then the Hessian matrix of $\Lambda + R_0$ is positive definite everywhere. A (unique) minimizer of $\Lambda + R_0$ fails to exist if, and only if, there exist vectors $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_L \in \mathbb{R}^{d'}$ such that $\boldsymbol{\theta}_y \neq \boldsymbol{\theta}_z$ for some class labels $1 \leq y < z \leq L$ and*

$$\boldsymbol{\theta}_{Y_i}^\top \mathbf{V}_i = \max_{y=1, \dots, L} \boldsymbol{\theta}_y^\top \mathbf{V}_i \text{ for } 1 \leq i \leq n. \quad (1.29)$$

Theorem 1.16. *Let $R = R_0 + R_1$ or $R = R_0 + R_2$. If $\sigma_1 > 0$ and $\tau_j > 0$ for $2 \leq j \leq d'$, then the Hessian matrix of $\Lambda + R$ is positive definite everywhere, and there exists a (unique) minimizer of $\Lambda + R$.*

PROOF OF THEOREM 1.15. The Hessian of $\Lambda + R_0$ at $\boldsymbol{\theta}$ is the sum of the positive semidefinite matrices $\mathbf{H}(\boldsymbol{\theta})$ and $\mathbf{H}_{R,0}(\boldsymbol{\theta})$. According to Theorem 1.10, for any $\boldsymbol{\delta} \in \mathbb{R}^{Ld'}$, it follows from $\boldsymbol{\delta}^\top \mathbf{H}(\boldsymbol{\theta})\boldsymbol{\delta} = 0$ that $\boldsymbol{\delta}_1 = \boldsymbol{\delta}_2 = \dots = \boldsymbol{\delta}_L$. But $\boldsymbol{\delta}^\top \mathbf{H}_{R,0}(\boldsymbol{\theta})\boldsymbol{\delta} = 2R_0(\boldsymbol{\delta}) = \sum_{j=1}^{d'} \sigma_j (\mathbf{1}_L^\top \boldsymbol{\delta}_{[j]})^2 = 0$ implies that $\mathbf{1}_L^\top \boldsymbol{\delta}_{[j]} = 0$ for $1 \leq j \leq d'$, so $\sum_{y=1}^L \boldsymbol{\delta}_y = L\boldsymbol{\delta}_z = \mathbf{0}$ for $1 \leq z \leq L$.

As to the existence of a unique minimizer, note that strict convexity of $f = R + R_0$ implies that it has either a unique minimizer or no minimizer at all. Hence existence of a minimizer is equivalent to (1.28). Suppose that the latter condition is violated, i.e. $\lim_{t \rightarrow \infty} \boldsymbol{\theta}^\top \nabla f(t\boldsymbol{\theta}) \leq 0$ for a fixed nonzero $\boldsymbol{\theta}$. Note that $\boldsymbol{\theta}^\top \nabla f(t\boldsymbol{\theta})$ is the sum of $\boldsymbol{\theta}^\top \mathbf{G}(t\boldsymbol{\theta})$ and $\boldsymbol{\theta}^\top \mathbf{G}_{R,0}(t\boldsymbol{\theta}) = 2tR_0(\boldsymbol{\theta})$. Moreover,

$$\begin{aligned} \boldsymbol{\theta}^\top \mathbf{G}(t\boldsymbol{\theta}) &= \sum_{i=1}^n \left(\sum_{y=1}^L \frac{\exp(t\boldsymbol{\theta}_y^\top \mathbf{V}_i)}{\sum_{z=1}^L \exp(t\boldsymbol{\theta}_z^\top \mathbf{V}_i)} \boldsymbol{\theta}_y^\top \mathbf{V}_i - \boldsymbol{\theta}_{Y_i}^\top \mathbf{V}_i \right) \\ &\rightarrow \sum_{i=1}^n \left(\max_{y=1, \dots, L} \boldsymbol{\theta}_y^\top \mathbf{V}_i - \boldsymbol{\theta}_{Y_i}^\top \mathbf{V}_i \right), \end{aligned}$$

1. Classifiers and P-Values

which is certainly nonnegative. Hence, our assumption entails that $R_0(\boldsymbol{\theta}) = 0$, i.e. $\sum_{y=1}^L \boldsymbol{\theta}_y = \mathbf{0}$ and (1.29). Since $\boldsymbol{\theta} \neq \mathbf{0}$ and $\sum_{y=1}^L \boldsymbol{\theta}_y = \mathbf{0}$, the subvectors $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_L$ cannot be all identical.

On the other hand, suppose that (1.29) holds for some $\boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_L^\top)^\top$ such that $\boldsymbol{\theta}_y \neq \boldsymbol{\theta}_z$ for some $1 \leq y < z \leq L$. These properties remain unchanged if we subtract $L^{-1} \sum_{y=1}^L \boldsymbol{\theta}_y$ from all subvectors $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_L$. But then we have a nonzero vector $\boldsymbol{\theta} \in \mathbb{R}^{Ld'}$ such that $\boldsymbol{\theta}^\top \mathbf{G}_{R,0}(t\boldsymbol{\theta}) = 0$ and $\boldsymbol{\theta}^\top \mathbf{G}(t\boldsymbol{\theta}) \rightarrow 0$ as $t \rightarrow \infty$. \square

PROOF OF THEOREM 1.16. Suppose that for some $\boldsymbol{\theta}, \boldsymbol{\delta} \in \mathbb{R}^{Ld'}$,

$$\boldsymbol{\delta}^\top \mathbf{H}(\boldsymbol{\theta}) \boldsymbol{\delta} = \boldsymbol{\delta}^\top \mathbf{H}_{R,0}(\boldsymbol{\theta}) \boldsymbol{\delta} = \boldsymbol{\delta}^\top \mathbf{H}_{R,k}(\boldsymbol{\theta}) \boldsymbol{\delta} = 0,$$

where $k = 1$ or $k = 2$. It follows from $\boldsymbol{\delta}^\top \mathbf{H}_{R,k}(\boldsymbol{\theta}) \boldsymbol{\delta} = 0$ that $\delta_{j,y} = 0$ for $2 \leq j \leq d'$ and $1 \leq y \leq L$. But then $\boldsymbol{\delta}^\top \mathbf{H}(\boldsymbol{\theta}) \boldsymbol{\delta} = 0$ is equivalent to $\delta_{1,1} = \delta_{1,2} = \dots = \delta_{1,L}$, because all vectors \mathbf{V}_i have first component one. Hence it follows from $\boldsymbol{\delta}^\top \mathbf{H}_{R,0}(\boldsymbol{\theta}) \boldsymbol{\delta} = \sigma_1(\mathbf{1}_L^\top \boldsymbol{\delta}_{[1]})^2 = 0$ that $\mathbf{1}_L^\top \boldsymbol{\delta}_{[1]} = L\delta_{1,y} = 0$ for $1 \leq y \leq L$.

For the existence of a unique minimizer we employ (1.28) again. \square

1.8.4. Some Comments on the Implementation in pvclass

Representations with matrices. For various reasons it is better to work with the parameter matrix $\underline{\boldsymbol{\theta}} = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_L] \in \mathbb{R}^{d' \times L}$. Then $\boldsymbol{\theta} = \text{vec}(\underline{\boldsymbol{\theta}})$. In R the operator $\text{vec}(\cdot)$ is implemented as `as.vector(\cdot)`. Together with the augmented data matrix $\underline{\mathbf{V}} = [\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_n]^\top$, one may write

$$\underline{\mathbf{F}} := [\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_1), \mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_2), \dots, \mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_n)]^\top = \underline{\mathbf{V}} \underline{\boldsymbol{\theta}},$$

and this simplifies various computations, e.g. the determination of

$$\underline{\mathbf{P}} := [\mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_1)), \mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_2)), \dots, \mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_n))]^\top,$$

considerably. Another trick to speed up computations is the well-known relation

$$\mathbf{a} \otimes \mathbf{b} = \text{vec}(\mathbf{b}\mathbf{a}^\top)$$

for arbitrary (column) vectors \mathbf{a} and \mathbf{b} . In particular, the gradient $\mathbf{G}(\boldsymbol{\theta})$ of Λ at $\boldsymbol{\theta}$ may be represented as

$$\begin{aligned}\mathbf{G}(\boldsymbol{\theta}) &= \sum_{i=1}^n (\mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) - \mathbf{e}_{L, Y_i}) \otimes \mathbf{V}_i \\ &= \text{vec} \left(\sum_{i=1}^n \mathbf{V}_i (\mathbf{p}(\mathbf{f}(\boldsymbol{\theta}, \mathbf{V}_i)) - \mathbf{e}_{L, Y_i})^\top \right) \\ &= \text{vec}(\mathbf{V}^\top (\underline{\mathbf{P}} - \underline{\mathbf{E}})),\end{aligned}$$

where $\underline{\mathbf{E}} := [\mathbf{e}_{L, Y_1}, \mathbf{e}_{L, Y_2}, \dots, \mathbf{e}_{L, Y_n}]^\top$. For the Hessian $\mathbf{H}(\boldsymbol{\theta})$ we also avoid the summation of n Kronecker products as follows: For $y, z \in \{1, 2, \dots, L\}$,

$$(\mathbf{H}(\boldsymbol{\theta}))_{ij} = \mathbf{V}^\top \left(((\delta_{yz} \mathbf{P}_y - \mathbf{P}_y \odot \mathbf{P}_z) \mathbf{1}_{d'}^\top) \odot \mathbf{V} \right),$$

where \odot stands for componentwise multiplication, and $\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_L$ are the columns of $\underline{\mathbf{P}}$.

Normalization of $\underline{\mathbf{F}}$. Having computed $\underline{\mathbf{F}} = (F_{iy})_{i \leq n, y \leq L}$, one can write $\underline{\mathbf{P}} = (P_{iy})_{i \leq n, y \leq L}$ as

$$P_{iy} = \exp(F_{iy}) / \sum_{z=1}^L \exp(F_{iz}),$$

and

$$\Lambda(\boldsymbol{\theta}) = \sum_{i=1}^n \left(\log \left(\sum_{z=1}^L \exp(F_{iz}) \right) - F_{i, Y_i} \right).$$

These representations become problematic numerically if some components of $\underline{\mathbf{F}}$ become very large or if a whole row of $\underline{\mathbf{F}}$ consists of very small (negative) values. To circumvent these problems, we subtract from each row of $\underline{\mathbf{F}}$ its maximum. This does not affect the previous expressions for P_{iy} or $\Lambda(\boldsymbol{\theta})$.

Choice of σ and τ . In our implementations in `pvc` class, we use three versions of penalized logistic regression, specified by the parameters `pen.method` and τ_o :

pen.method	R	σ	τ
vectors	$R_0 + R_1$	$\mathbf{1}_{d'}$	$(\tau_o S_j)_{j=1}^{d'}$
simple	$R_0 + R_2$	$\mathbf{e}_{d', 1}$	$(\tau_o S_j)_{j=1}^{d'}$
none	R_0	$\mathbf{1}_{d'}$	—

Here S_j is the sample standard deviation (within groups) of the j -th components of the vectors \mathbf{V}_i .

1. Classifiers and P-Values

Starting values. In the standard model with $\mathbb{P}(Y = y) = w_y$ and $\mathcal{L}(X | Y = y) = \mathcal{N}_d(\boldsymbol{\mu}_y, \boldsymbol{\Sigma})$,

$$\mathbb{P}(Y = y | \mathbf{X} = \mathbf{x}) = \frac{\exp(\log w_y - \boldsymbol{\mu}_y^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_y / 2 + \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_y)}{\sum_{z=1}^L \exp(\log w_z - \boldsymbol{\mu}_z^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_z / 2 + \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_z)}.$$

Hence a possible starting point for iterative minimization algorithms is given by

$$\boldsymbol{\theta}_y^{(0)} = \boldsymbol{\theta}_y^{(*)} - \mathbf{M}^{(*)} \quad \text{with} \quad \boldsymbol{\theta}_y^{(*)} := \begin{pmatrix} \log \hat{w}_y - \hat{\boldsymbol{\mu}}_y^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_y / 2 \\ \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_y \end{pmatrix},$$

where \hat{w}_y , $\hat{\boldsymbol{\mu}}_y$ and $\hat{\boldsymbol{\Sigma}}$ are the usual parameter estimators in linear discriminant analysis, while $\mathbf{M}^{(*)} \in \mathbb{R}^{d'}$ is a centering vector depending on the type of regularization. Its first component equals the average of the components of $\boldsymbol{\theta}_{[1]}^{(*)}$. For $2 \leq j \leq d'$, the j -th component of $\mathbf{M}^{(*)}$ is either the mean (`pen.method = vectors`) or the median (`pen.method = simple`) of the components of $\boldsymbol{\theta}_{[j]}^{(*)}$.

Solution paths. Suppose that f_1, f_2 are two convex and twice continuously differentiable functionals on \mathbb{R}^q . Suppose further that for any $t > 0$, the functional $f_1 + t f_2$ is coercive with positive definite Hessian matrix $D^2 f_1 + t D^2 f_2$ everywhere. This entails that for each $t > 0$ there exists a unique minimizer $\boldsymbol{\theta}(t) \in \mathbb{R}^q$ of $f_1 + t f_2$; indeed $\boldsymbol{\theta}(t)$ is the unique solution $\boldsymbol{\theta}$ of the equation

$$\nabla f_1(\boldsymbol{\theta}) + t \nabla f_2(\boldsymbol{\theta}) = 0.$$

It follows from the implicit mapping theorem, applied to the function $\mathbb{R}^d \times \mathbb{R} \ni (\mathbf{x}, t) \mapsto \nabla f_1(\mathbf{x}) + t \nabla f_2(\mathbf{x}) \in \mathbb{R}^d$, that $(0, \infty) \ni t \mapsto \boldsymbol{\theta}(t)$ is also continuously differentiable with derivative

$$\boldsymbol{\theta}'(t) = - (D^2 f_1(\boldsymbol{\theta}(t)) + t D^2 f_2(\boldsymbol{\theta}(t)))^{-1} \nabla f_2(\boldsymbol{\theta}(t)).$$

These considerations are useful when minimizing $\Lambda + R_0 + \tau_o R_k$ for different values of $\tau_o > 0$. Having determined the minimizer $\boldsymbol{\theta}(\tau_o)$ for some value of τ_o , a good starting point for the Newton procedure with τ_* close to τ_o is given by

$$\begin{aligned} \boldsymbol{\theta}^{(0)} &:= \boldsymbol{\theta}(\tau_o) - (\tau_* - \tau_o) (\mathbf{H}(\boldsymbol{\theta}(\tau_o)) + \mathbf{H}_{R,0}(\boldsymbol{\theta}(\tau_o)) \\ &\quad + \tau_o \mathbf{H}_{R,k}(\boldsymbol{\theta}(\tau_o)))^{-1} \mathbf{G}_{R,k}(\boldsymbol{\theta}(\tau_o)) \\ &= \boldsymbol{\theta}(\tau_o) - (\tau_*/\tau_o - 1) (\mathbf{H}(\boldsymbol{\theta}(\tau_o)) + \mathbf{H}_{R,0}(\boldsymbol{\theta}(\tau_o)) \\ &\quad + \tau_o \mathbf{H}_{R,k}(\boldsymbol{\theta}(\tau_o)))^{-1} (\tau_o \mathbf{G}_{R,k}(\boldsymbol{\theta}(\tau_o))). \end{aligned}$$

2. Choice of Tuning Parameters

Some of the test statistics we use depend on a tuning parameter such as the k in the nearest neighbor method or the penalty parameter τ in the logistic regression. We want to choose them in a data-driven way which preserves the symmetry in $(\mathbf{X}_{I(j)})_{j=1}^{N_\theta}$.

Our first approach was to optimize the estimated expectations of the p-values. To choose the tuning parameter for the p-value $\pi_\theta(\mathbf{X}, \mathcal{D})$, we add the new observation \mathbf{X} to the training data with class label θ . Then we search for the parameter which minimizes the sum of the cross-validated p-values

$$\sum_{i=1}^n \pi_\theta(\mathbf{X}_i, \mathcal{D}).$$

Unfortunately, this method chooses mostly small values for k or the regularization parameter. The reason for this is overfitting. It selects the parameter for which the classes are separated best. However, it is not taken into account how the p-values change, if the training data vary slightly. For example if we add a penalty term in the logistic regression, the separation of the training data gets worse, but we gain stability.

2.1. Stability

In a second approach we want to maximize the stability, i.e. $T_\theta(\mathbf{X}, \mathcal{D})$ should take big values for observations not belonging to class θ . To find the tuning parameter which maximizes the stability, we add the new observation \mathbf{X} to the training data with class label θ . Then we compute for all training observations with $Y_i \neq \theta$ the test statistic

$$T_\theta^{(\tau)}(\mathbf{X}_i, \mathcal{D}_i(\mathbf{X}_i, \mathbf{X}, \theta)),$$

where $\mathcal{D}_i(\mathbf{X}_i, \mathbf{X}, \theta)$ denotes the training data after adding the observation (\mathbf{X}, θ) and setting the class label of observation \mathbf{X}_i to θ . Then we take the sum of these values

$$S(\tau, \mathbf{X}, \theta) := \sum_{i: Y_i \neq \theta} T_\theta^{(\tau)}(\mathbf{X}_i, \mathcal{D}_i(\mathbf{X}_i, \mathbf{X}, \theta))$$

2. Choice of Tuning Parameters

and search for the parameter τ^* which maximizes $S(\tau, \mathbf{X}, \theta)$:

$$\tau^*(\mathbf{X}, \theta) := \arg \max_{\tau} S(\tau, \mathbf{X}, \theta).$$

2.1.1. Subsampling

To determine the optimal tuning parameters for a new observation \mathbf{X} and all potential class memberships, the test statistic has to be computed $(L-1) \cdot n \cdot l$ times, where l is the number of tuning parameters from which we want to choose the optimal one. This can be very computer-intensive, especially for penalized logistic regression in high dimensions. One way to reduce computation time is to take a random subset of the training data containing m observations per class $b \neq \theta$ and compute $T_{\theta}^{(\tau)}(\mathbf{X}_i, \mathcal{D}_i(\mathbf{X}_i, \mathbf{X}, \theta))$ only for these observations. Subsampling is particularly useful for large training sets.

2.1.2. Extended Golden Section Search

Another way to reduce computation time is to consider only few values for τ instead of a whole grid. We observed in simulated and real data examples of penalized logistic regression that $S(\tau, \mathbf{X}, \theta)$ is a unimodal function of τ for fixed \mathbf{X} and θ , at least for reasonable intervals $[a, b]$.

Now suppose that $\tau^* \in [a, b]$ and $S(\tau): [a, b] \rightarrow \mathbb{R}$ is *unimodal*, i.e. S is strictly increasing on $[a, \tau^*]$ and strictly decreasing on $[\tau^*, b]$. In this case we can use the golden section search.

Let $a \leq q < r < s < t \leq b$ such that $\tau^* \in [q, t]$. At the beginning of the algorithm we set $q = a$ and $t = b$. Then $S(r) \leq S(s)$ implies $\tau^* \in [r, t]$ and we replace $(q < r < s < t)$ by $(r < s < s' < t)$. In case of $S(r) > S(s)$ we know that $\tau^* \in [q, s]$ and replace $(q < r < s < t)$ by $(q < r' < r < s)$.

The new point r' or s' , respectively, can be chosen in different ways. For the golden section search (Kiefer, 1953) we choose $r' = Cq + Bs$ or $s' = Br + Ct$, respectively, with $C = (\sqrt{5} - 1)/2 \approx 0.618$ and $B = 1 - C \approx 0.382$, i.e. we divide the interval in the golden ratio. In this way we can guarantee that the quadruples $(q < r < s < t)$, $(r < s < s' < t)$ and $(q < r' < r < s)$ differ only by an affine transformation. Moreover, we get in each step a reduction of the interval length by the factor C . We stop the algorithm when $t - q \leq \delta$ and end up with an interval containing τ^* with length smaller than $C\delta$. Now we set $\tau^* = \arg \max(S(q), S(r), S(s), S(t))$. With this algorithm we don't get the exact argument of the maximum of S . However, the resulting p-values do not depend too sensitively on the exact choice of τ and it suffices to choose a τ which is not too far from the arg max. For the same reason, we can choose a rather large value for δ .

Since we don't know an interval $[a, b]$ which contains τ^* , we extend the golden section search. First we search for such an interval and then we start the

golden section search. As lower endpoint a we could take 0. However, to avoid numerical problems we suggest to take a small value greater than 0, e.g. 1. But a should be small enough, such that we may assume that $\tau^* \geq a$. To find the upper endpoint, we start with some point b and divide the interval $[a, b]$ in the golden ratio, i.e. $s = Ba + Cb$. If $S(s) \geq S(b)$, $\tau^* \in [a, b]$ and we can start the golden section search. Otherwise, τ^* could be greater than b and we define a new upper endpoint $b' = b + C(b - a)$. Then we iterate these steps until $S(s) \geq S(b)$. Pseudocode for the resulting algorithm is given in Algorithm 1. In the last step we assign $\tau^* = \arg \max(S(q), S_r, S_s)$, because $S(t) \leq S_s$ by construction of the algorithm.

Algorithm 1 $\tau^* \leftarrow \text{extendedGoldenSection}(S, a, b, \delta)$

```

 $(B, C) \leftarrow (0.382, 0.618)$ 
 $s \leftarrow Ba + Cb$ 
 $(S_s, S_b) \leftarrow (S(s), S(b))$ 
if  $S_s < S_b$  then
     $\tau^* \leftarrow \text{extendedGoldenSection}(S, a, b + C(b - a), \delta)$ 
    return  $\tau^*$ 
end if
 $(q, r, t) \leftarrow (a, Ca + Bb, b)$ 
 $S_r \leftarrow S(r)$ 
while  $t - q > \delta$  do
    if  $S_r > S_s$  then
         $(r, s, t) \leftarrow (Cq + Bs, r, s)$ 
         $(S_r, S_s) \leftarrow (S(r), S_r)$ 
    else
         $(q, r, s) \leftarrow (r, s, Br + Ct)$ 
         $(S_r, S_s) \leftarrow (S_s, S(s))$ 
    end if
end while
 $\tau^* \leftarrow \arg \max(S(q), S_r, S_s)$ 
return  $\tau^*$ 

```

2.2. Dimension Reduction

In very high-dimensional settings the computation of the p-values and the choice of τ , even with subsampling, become too computer-intensive. Therefore we reduce the dimension of the data before applying the procedure. The aim of this reduction is not to eliminate all noise variables, but to find a small enough subset of predictors for which our procedure is computationally feasible. We use l_1 -penalized multicategory logistic regression to determine the

2. Choice of Tuning Parameters

subset.

To preserve the symmetry in $(\mathbf{X}_{I(j)})_{j=1}^{N_\theta}$ we have two options. We can split the training set and use one part for reduction and the other as training set. If the training set is not big enough, we add the new observation to the training set and use the whole set for the reduction. The drawback of the latter method is that we have to do the reduction separately for each new observation and potential class membership.

2.3. Numerical Examples

2.3.1. Simulated Data

Example 2.1. Consider $L = 2$ classes with $P_\theta = \mathcal{N}_{100}(\boldsymbol{\mu}_\theta, \mathbb{I}_{100})$, where $\boldsymbol{\mu}_1 = (1, 0.5, 0.25, 0.125, 0, \dots, 0)^\top$ and $\boldsymbol{\mu}_2 = -\boldsymbol{\mu}_1$. We simulated $N_1 = N_2 = 100$ training observations per class and computed $\pi_\theta^{(\tau)}(\mathbf{X}, \mathcal{D}(\mathbf{X}, \theta))$ for 100 test observations per class and $\tau = 1, 2, \dots, 100$. Here $\pi_\theta^{(\tau)}$ denotes the p-value based on penalized logistic regression as described in Section 1.4. Figure 2.1 shows the rate of uniquely correct classified test observations for five different training and test sets. The missclassification rates depend heavily on the training set. Therefore we simulated 100 training and test sets and averaged the missclassification rates. The result is also shown in Figure 2.1 (bottom right). The corresponding distributions of τ^* are drawn as bar plots. Note that they are scaled to 1 and do not correspond to the scale of the y-axis. To determine the distribution of τ^* we used 20 training sets and 100 test observations per class and training set.

Example 2.2. Next we consider an example with $L = 5$ classes, $P_\theta = \mathcal{N}_{100}(\boldsymbol{\mu}_\theta, \mathbb{I}_{100})$ and $\boldsymbol{\mu}_\theta = (\mathbf{0}_{4(\theta-1)}, 2, 1, 0.5, 0.25, 0, \dots, 0)^\top$, where $\mathbf{0}_j$ denotes the row vector with j zeros. We simulated $N_\theta = 100$ training observations and 40 test observations per class. Figure 2.2 shows the rate of uniquely correct classified test observations for five different training and test sets and averaged over 100 training and test sets (bottom right). The corresponding distributions of τ^* are drawn as bar plots. To determine the distribution of τ^* we used 5 training sets and 40 test observations per class and training set.

The choice of τ is good in both examples, but the procedure tends to pick slightly too large values for τ , i.e. to regularize slightly too strong. This is better than too little regularization, since a strong regularization avoids overfitting, which is a big problem, especially in high-dimensional settings.

Subsampling increases the variability of τ^* . But the results remain quite good for reasonable choices of m , e.g. $m = 10$.

Note that we averaged the rates of uniquely correct classified test observations over all classes, since these examples are symmetric.

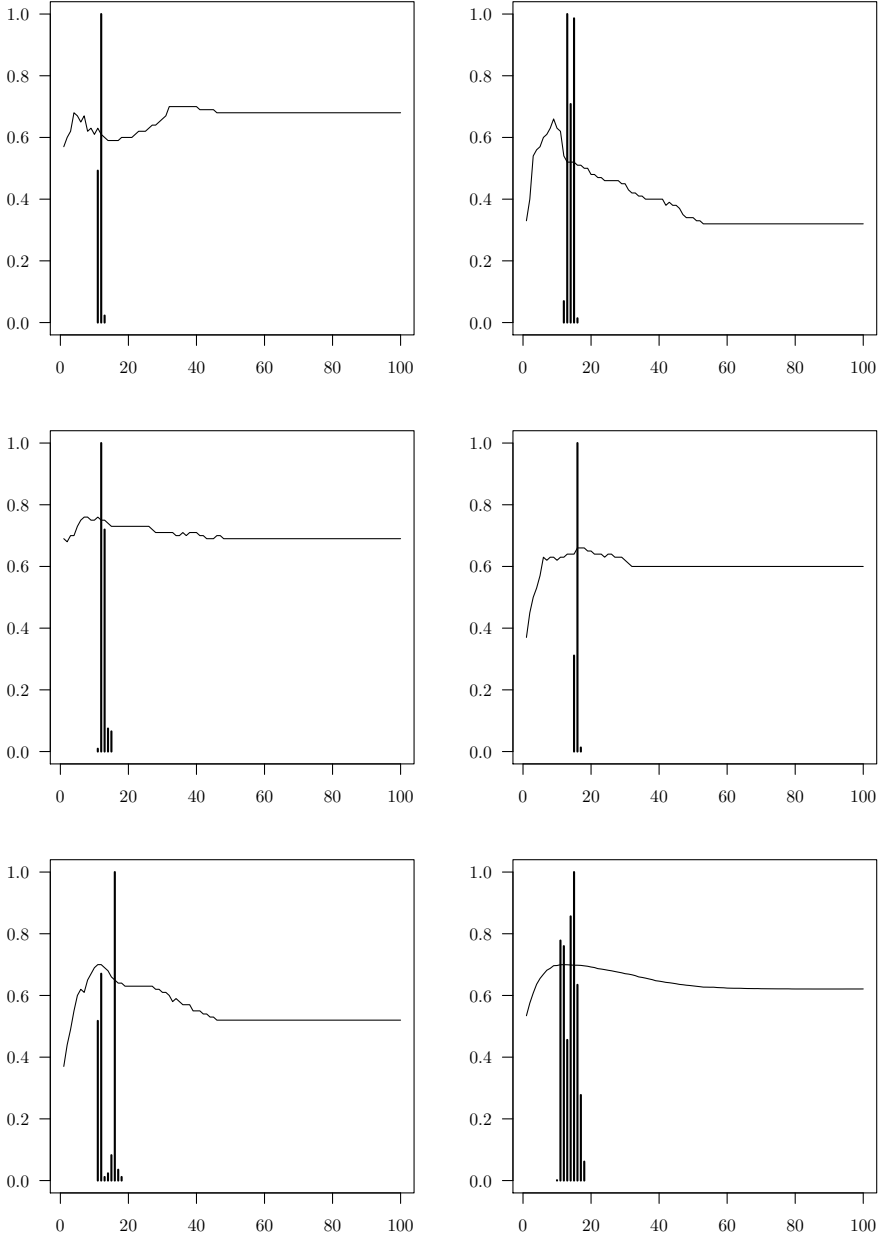


Figure 2.1.: Rates of uniquely correct classified test observations and distributions of τ^* for Example 2.1.

2. Choice of Tuning Parameters

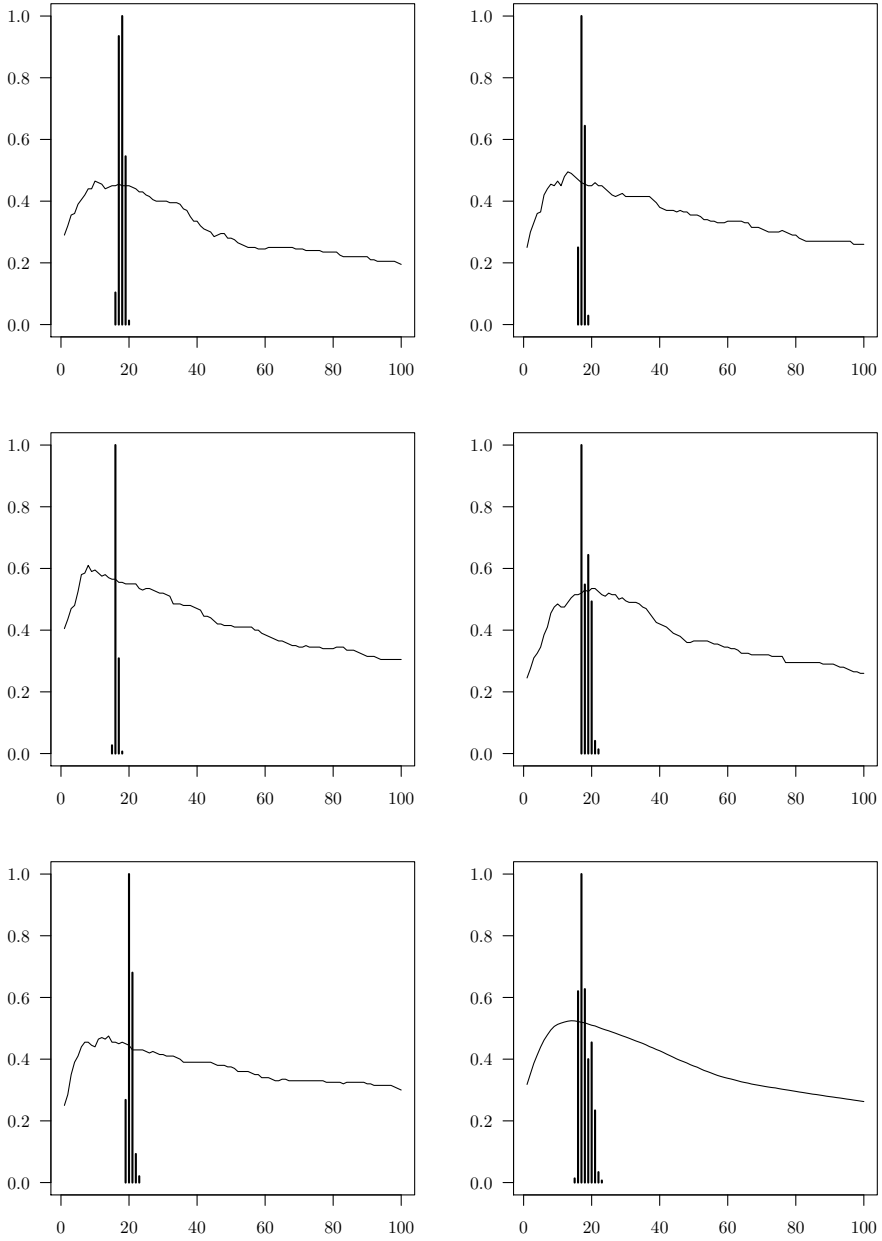


Figure 2.2.: Rates of uniquely correct classified test observations and distributions of τ^* for Example 2.2.

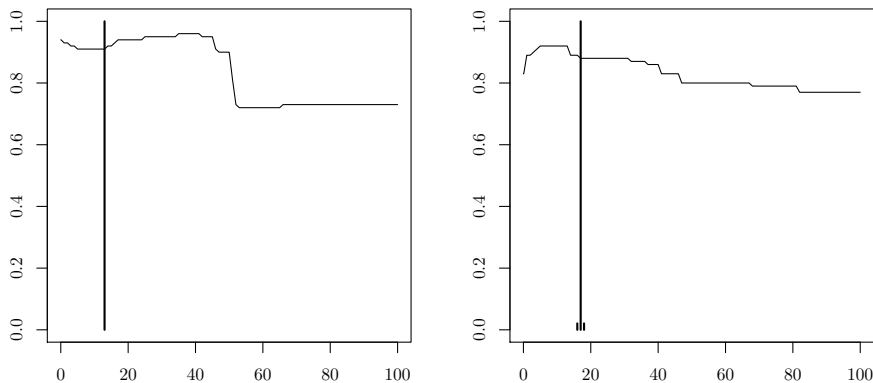


Figure 2.3.: Rates of uniquely correct classified test observations and distribution of τ^* for Example 2.3.

2.3.2. Real Data

Example 2.3 (Internet Ad). We consider data of images on Internet pages provided by Kushmerick (1999). The binary response indicates whether the image is an advertisement. The preprocessed data set (Friedman et al., 2010) without missing values contains 1978 “nonads” and 381 ads. Three of the 1430 features are continuous, the others are binary. Computing p-values for such a high-dimensional problem would be too computer-intensive. Therefore we reduce the dimension as described in Section 2.2 using the R-package glmnet (Friedman et al., 2010). We split the data in a set for the reduction with 500 “nonads” and 100 ads, a test set with 100 observations per class and a training set with 1378 “nonads” and 181 ads.

For the reduction, we choose a penalty parameter for which we end up with a problem of dimension 120. Figure 2.3 shows the rate of uniquely correct classified test observations of class “nonads” (left) and ads (right). The distributions of τ^* are drawn as bar plots. They were determined using subsampling with $m = 100$.

Again, τ^* is not at the arg max of the rate of unique correct classification, but the rate for τ^* is near to the maximum. For class ads the procedure again regularizes slightly too strong.

Example 2.4 (Mushrooms). The UCI Machine Learning Repository (Bache and Lichman, 2013) provides data of hypothetical samples corresponding to 23 species of gilled mushrooms. The binary response indicates whether the species is edible or poisonous. Most of the 22 categorical features describe the shape or the color. We removed 2 features. One has missing values and the other one takes only one value. After creating dummy-variables we end up

2. Choice of Tuning Parameters

b	$P(b, \{\})$	$P(b, \{1\})$	$P(b, \{2\})$	$P(b, \{1, 2\})$
1	0.07	0.93	0.00	0
2	0.08	0.00	0.92	0

Table 2.1.: Missclassification rates for $\tau = 40$ in Example 2.4.

with a problem of dimension 91. The data set contains observations of 4208 edible and 3916 poisonous mushrooms. We picked randomly 100 observations per class for the test set and used the rest as training set.

We computed the p-values for $\tau = 0.001, 0.01, 0.1, 1, 2, 3, \dots, 100$. Table 2.1 shows the missclassification rates for $\tau = 40$. The choice of τ has no big influence on the result in this example. The rate of uniquely correct classified observations varies only between 0.9 and 0.94. All the p-values for the wrong classes are smaller than $\alpha = 0.05$. The two classes are perfectly separated, but for some of the observations both null hypotheses are rejected. They are located between the two class centers at a big distance from most of the training observations. Therefore both possible class memberships seem implausible for these observations.

Example 2.5 (Buerk). For the hospital data described in Section 1.7 we computed p-values for $\tau = 1, 2, \dots, 100$. The amount of regularization has no influence on the missclassification rates in this example.

3. Central Limit Theorems

In this chapter we derive two central limit theorems. First, we consider linear discriminant analysis and describe the asymptotic distribution of missclassification rates and cross-validated estimators thereof. Second, we consider p-values based on the plug-in statistic for the standard model and prove a central limit theorem for conditional inclusion probabilities and empirical conditional inclusion probabilities, which can be interpreted as estimators of the former.

We consider $L = 2$ classes with distributions P_1 and P_2 on $\mathcal{X} = \mathbb{R}^d$, which differ only by a shift. For notational convenience we assume without loss of generality that $Y_1 = 1$ and $Y_2 = 2$. Let $\mathbb{E}(\mathbf{X}_1) = \boldsymbol{\mu}_1$ and $\mathbb{E}(\mathbf{X}_2) = \boldsymbol{\mu}_2$ denote the mean vectors and suppose that

$$\boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2.$$

Let $\tilde{\mathbf{X}}_i := \mathbf{X}_i - \boldsymbol{\mu}_{Y_i}$ denote the centered observation and $\boldsymbol{\Sigma} = \text{Var}(\mathbf{X}_1) = \text{Var}(\mathbf{X}_2)$ the common positive definite covariance matrix.

Moreover, let $P_0 := \mathcal{L}(\boldsymbol{\Sigma}^{-1/2}(\mathbf{X}_1 - \boldsymbol{\mu}_1)) = \mathcal{L}(\boldsymbol{\Sigma}^{-1/2}(\mathbf{X}_2 - \boldsymbol{\mu}_2)) = \mathcal{L}(\mathbf{Z})$ with $\mathbf{Z} := \boldsymbol{\Sigma}^{-1/2}(\mathbf{X}_1 - \boldsymbol{\mu}_1)$ and its first component Z_1 .

We assume that \mathbf{Z} has a differentiable Lebesgue density $f_{\mathbf{Z}} > 0$ and for all $\mathbf{U} \in \mathbb{R}^{d \times (d-1)}$ with $\mathbf{U}^\top \mathbf{U} = \mathbb{I}_{d-1}$

$$K(\mathbf{U}) := \int_{\mathbb{R}^{d-1}} \sup_{\mathbf{a} \in \mathbf{U}^\perp} \|\nabla f_{\mathbf{Z}}(\mathbf{a} + \mathbf{U}\mathbf{z})\| (1 + \|\mathbf{z}\|)^2 \, d\mathbf{z} < \infty. \quad (3.1)$$

Here \mathbb{I}_d denotes the d -dimensional identity matrix and \mathbf{U}^\perp the orthogonal complement of the column space of \mathbf{U} . We define $\boldsymbol{\beta} := \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$, such that $\|\boldsymbol{\beta}\| = D_{\boldsymbol{\Sigma}}(\boldsymbol{\mu}_2, \boldsymbol{\mu}_1)$. Except for the proof of Lemma 3.13. There we define $\boldsymbol{\beta} := \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta)$.

In this chapter we consider the class labels Y_1, Y_2, \dots, Y_n as fixed while $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ and (\mathbf{X}, Y) are independent. We use $\hat{w}_\theta := N_\theta/n$ as estimator for the prior weights w_θ . Note that the choice of \hat{w}_θ is relevant only for the linear discriminant analysis. The p-values based on the plug-in statistic for the standard Gaussian model for two classes do not depend on the prior weights w_θ . Asymptotic statements are meant as

$$n \rightarrow \infty \quad \text{and} \quad \hat{w}_\theta \rightarrow w_\theta \quad \text{for all } \theta \in \mathcal{Y}.$$

3. Central Limit Theorems

Convergence in probability is denoted by \rightarrow_p , convergence in law by $\rightarrow_{\mathcal{L}}$ and almost sure convergence by $\rightarrow_{\text{a.s.}}$. Generally we denote with $\widehat{\mathbf{v}}$ an estimator of the vector $\mathbf{v} = (v_1, v_2, \dots, v_d)^\top$, by $\Delta_{\mathbf{v}}$ the scaled difference $\sqrt{n}(\widehat{\mathbf{v}} - \mathbf{v}) = (\Delta_{\mathbf{v},1}, \Delta_{\mathbf{v},2}, \dots, \Delta_{\mathbf{v},d})^\top$ and $\mathbf{Y}_{n,i}^{\mathbf{v}}$ is a summand depending on \mathbf{X}_i such that $\Delta_{\mathbf{v}} = \sum_{i=1}^n \mathbf{Y}_{n,i}^{\mathbf{v}} + o_p(1)$. Similarly, $\widehat{\mathbf{A}}$ is an estimator for the matrix

$$\mathbf{A} = \begin{pmatrix} A_{1,1} & \cdots & A_{1,d} \\ \vdots & \ddots & \vdots \\ A_{d,1} & \cdots & A_{d,d} \end{pmatrix},$$

$\Delta_{\mathbf{A}}$ the scaled difference $\sqrt{n}(\widehat{\mathbf{A}} - \mathbf{A})$ and $\mathbf{Y}_{n,i}^{\mathbf{A}}$ is a summand depending on \mathbf{X}_i such that $\Delta_{\mathbf{A}} = \sum_{i=1}^n \mathbf{Y}_{n,i}^{\mathbf{A}} + o_p(1)$. Moreover, $\mathbf{v}_{i:j} = (v_i, v_{i+1}, \dots, v_j)^\top$ is a vector consisting of the components i to j of \mathbf{v} .

The density of a random variable ξ is denoted by f_ξ . If ξ is one-dimensional F_ξ denotes its distribution function.

For random matrices $\widehat{\Sigma}$, $\text{Var}(\widehat{\Sigma}) := \text{Var}(\text{vec}(\widehat{\Sigma}))$, where $\text{vec}(\mathbf{M})$ denotes a vector which is formed by stacking the columns of a matrix \mathbf{M} (from left to right). To formulate some of the results we use the Kronecker product \otimes defined on page 29.

We denote the symmetric difference of two sets A and B with $A \triangle B := (A \setminus B) \cup (B \setminus A)$.

3.1. Half-Spaces

Let \mathcal{H} denote the collection of all half-spaces in \mathbb{R}^d . A half-space is a set of the form

$$H(\boldsymbol{\beta}, \gamma) := \{\mathbf{z} \in \mathbb{R}^d : \boldsymbol{\beta}^\top \mathbf{z} + \gamma \leq 0\}$$

for $\boldsymbol{\beta} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ and $\gamma \in \mathbb{R}$.

The missclassification rates of the standard linear classifier can be written as a probability measure of a random subspace. Therefore we need some results about random half-spaces to describe asymptotic properties of missclassification rates. These results are also useful to describe the asymptotic behavior of inclusion probabilities for the p-values based on the plug-in statistic for the standard model.

3.1.1. Root- n -Consistency

The following lemma shows the root- n -consistency of the probability measure of random half-spaces under certain conditions.

Lemma 3.1. Let $a \in \mathbb{R}$, $\mathbf{b} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ and $\boldsymbol{\psi} := (\mathbf{b}^\top, a)^\top$. Further let \hat{a} and $\hat{\mathbf{b}}$ be random variables such that $\boldsymbol{\Delta}_\psi := \sqrt{n} \begin{pmatrix} \hat{\mathbf{b}} - \mathbf{b} \\ \hat{a} - a \end{pmatrix} = O_p(1)$. Additionally let \mathbf{X} be a random vector in \mathbb{R}^d with measure P and differentiable density f which satisfies

$$\|\mathbb{E}(\mathbf{X} \mid \mathbf{b}^\top \mathbf{X} = -a)\| < \infty \quad (3.2)$$

and for all $\mathbf{U} \in \mathbb{R}^{d \times (d-1)}$ with $\mathbf{U}^\top \mathbf{U} = \mathbb{I}_{d-1}$

$$K(\mathbf{U}) := \int_{\mathbb{R}^{d-1}} \sup_{\boldsymbol{\alpha} \in \mathbf{U}^\perp} \|\nabla f(\mathbf{a} + \mathbf{U}\mathbf{z})\| (1 + \|\mathbf{z}\|)^2 d\mathbf{z} < \infty. \quad (3.3)$$

Then

$$\sqrt{n}P(\{\mathbf{x} \in \mathbb{R}^d : \hat{\mathbf{b}}^\top \mathbf{x} + \hat{a} \leq 0\} \triangle \{\mathbf{x} \in \mathbb{R}^d : \mathbf{b}^\top \mathbf{x} + a \leq 0\}) = O_p(1) \quad (3.4)$$

and

$$\begin{aligned} & \sqrt{n}(P\{\mathbf{x} \in \mathbb{R}^d : \hat{\mathbf{b}}^\top \mathbf{x} + \hat{a} \leq 0\} - P\{\mathbf{x} \in \mathbb{R}^d : \mathbf{b}^\top \mathbf{x} + a \leq 0\}) \\ &= \mathbf{c}^\top \boldsymbol{\Delta}_\psi + O_p(n^{-1/2}) \end{aligned} \quad (3.5)$$

with $\mathbf{c}^\top = -f_{\mathbf{b}^\top \mathbf{X}}(-a)\mathbb{E}((\mathbf{X}^\top, 1) \mid \mathbf{b}^\top \mathbf{X} = -a)$, where $f_{\mathbf{b}^\top \mathbf{X}}$ denotes the density of $\mathbf{b}^\top \mathbf{X}$.

Remark 3.2. The constant a in the previous lemma can be replaced by a deterministic convergent sequence $a_n \rightarrow a$. We will use the lemma with a sequence depending on \hat{w}_1/\hat{w}_2 .

Remark 3.3. Condition (3.3) is satisfied for the multivariate t -Distribution with density

$$f(\mathbf{x}) = \det(\boldsymbol{\Sigma})^{-1/2} g_\nu((\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})),$$

for a mean vector $\boldsymbol{\mu} \in \mathbb{R}^d$, a nonsingular covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$, $\nu > 1$ and $g_\nu(s) := C_\nu^{-1}(1 + s/\nu)^{-(\nu+d)/2}$ with some normalizing constant $C_\nu > 0$.

PROOF OF REMARK 3.3. Without loss of generality let $\boldsymbol{\mu} = \mathbf{0}$ and $\boldsymbol{\Sigma} = \mathbb{I}_d$. Then

$$f(\mathbf{x}) = C_\nu^{-1}(1 + \nu^{-1}\|\mathbf{x}\|^2)^{-(\nu+d)/2}$$

and

$$\|\nabla f(\mathbf{x})\| = \frac{\nu + d}{C_\nu \nu} (1 + \nu^{-1}\|\mathbf{x}\|^2)^{-(\nu+d+2)/2} \|\mathbf{x}\|.$$

3. Central Limit Theorems

For $\mathbf{y} \in \arg \max_{\mathbf{x} \in \mathbb{R}^d} \|\nabla f(\mathbf{x})\|$ and $c := \|\mathbf{y}\|$,

$$\int_{\mathbb{R}^{d-1}} \sup_{\mathbf{a} \in \mathbf{U}^\perp} \|\nabla f(\mathbf{a} + \mathbf{U}\mathbf{z})\| (1 + \|\mathbf{z}\|)^2 \mathbb{1}\{\|\mathbf{z}\| \leq c\} \, d\mathbf{z} < \infty.$$

If $\|\mathbf{z}\| > c$, $\|\nabla f(\mathbf{a} + \mathbf{U}\mathbf{z})\| \leq \|\nabla f(\mathbf{U}\mathbf{z})\|$ for all $\mathbf{a} \in \mathbf{U}^\perp$. This implies for $c' > c$,

$$\begin{aligned} & \int_{\mathbb{R}^{d-1}} \sup_{\mathbf{a} \in \mathbf{U}^\perp} \|\nabla f(\mathbf{a} + \mathbf{U}\mathbf{z})\| (1 + \|\mathbf{z}\|)^2 \mathbb{1}\{c < \|\mathbf{z}\| \leq c'\} \, d\mathbf{z} \\ & \leq \int_{\mathbb{R}^{d-1}} \|\nabla f(\mathbf{U}\mathbf{z})\| (1 + \|\mathbf{z}\|)^2 \mathbb{1}\{c < \|\mathbf{z}\| \leq c'\} \, d\mathbf{z} \\ & = \frac{\nu + d}{C_\nu \nu} \int_{\mathbb{R}^{d-1}} (1 + \nu^{-1} \|\mathbf{z}\|^2)^{-(\nu+d+2)/2} \|\mathbf{z}\| (1 + \|\mathbf{z}\|)^2 \mathbb{1}\{c < \|\mathbf{z}\| \leq c'\} \, d\mathbf{z} \\ & = \frac{(d-1)\tau_{d-1}(\nu + d)}{C_\nu \nu} \int_{\mathbb{R}} \frac{r(1+r)^2 r^{d-2}}{(1 + \nu^{-1} r^2)^{(\nu+d+2)/2}} \mathbb{1}\{c < r \leq c'\} \, dr, \end{aligned}$$

where $\tau_{d-1} := \pi^{(d-1)/2} \Gamma((d+1)/2)^{-1}$ is the volume of the $(d-1)$ -dimensional unit sphere. By monotone convergence,

$$\begin{aligned} & \int_{\mathbb{R}^{d-1}} \sup_{\mathbf{a} \in \mathbf{U}^\perp} \|\nabla f(\mathbf{a} + \mathbf{U}\mathbf{z})\| (1 + \|\mathbf{z}\|)^2 \mathbb{1}\{c < \|\mathbf{z}\|\} \, d\mathbf{z} \\ & = \frac{(d-1)\tau_{d-1}(\nu + d)}{C_\nu \nu} \int_{\mathbb{R}} \frac{r(1+r)^2 r^{d-2}}{(1 + \nu^{-1} r^2)^{(\nu+d+2)/2}} \mathbb{1}\{c < r\} \, dr < \infty. \end{aligned}$$

□

PROOF OF LEMMA 3.1. Suppose that $\mathbf{b} = \mathbf{e}_1$, the first standard unit vector. Then

$$\begin{aligned} & P\{\mathbf{x} \in \mathbb{R}^d: \widehat{\mathbf{b}}^\top \mathbf{x} + \widehat{a} \leq 0\} - P\{\mathbf{x} \in \mathbb{R}^d: \mathbf{b}^\top \mathbf{x} + a \leq 0\} \\ & = P\left\{\left(\boldsymbol{\psi} + n^{-\frac{1}{2}} \boldsymbol{\Delta}_\psi\right)^\top \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \leq 0\right\} - P\left\{\boldsymbol{\psi}^\top \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \leq 0\right\} \\ & = P\left\{x_1 + a > 0 \geq x_1 + a + n^{-\frac{1}{2}} \boldsymbol{\Delta}_\psi^\top \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}\right\} \\ & \quad - P\left\{x_1 + a \leq 0 < x_1 + a + n^{-\frac{1}{2}} \boldsymbol{\Delta}_\psi^\top \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}\right\} \end{aligned}$$

$$\begin{aligned}
&= P\left\{x_1 + a > 0 \geq x_1(1 + n^{-\frac{1}{2}}\Delta_{\psi,1}) + a + n^{-\frac{1}{2}}\left(\Delta_{\psi,d+1} + \sum_{i=2}^d \Delta_{\psi,i}x_i\right)\right\} \\
&\quad - P\left\{x_1 + a \leq 0 < x_1(1 + n^{-\frac{1}{2}}\Delta_{\psi,1}) \right. \\
&\quad \left. + a + n^{-\frac{1}{2}}\left(\Delta_{\psi,d+1} + \sum_{i=2}^d \Delta_{\psi,i}x_i\right)\right\}.
\end{aligned}$$

With

$$\xi := \frac{-a - n^{-\frac{1}{2}}(\Delta_{\psi,d+1} + \sum_{i=2}^d \Delta_{\psi,i}x_i)}{1 + n^{-\frac{1}{2}}\Delta_{\psi,1}}$$

and the interval $I(\xi) := (\min(-a, \xi), \max(-a, \xi)]$ we may write

$$\begin{aligned}
&P\{\mathbf{x} \in \mathbb{R}^d: \widehat{\mathbf{b}}^\top \mathbf{x} + \widehat{a} \leq 0\} - P\{\mathbf{x} \in \mathbb{R}^d: \mathbf{b}^\top \mathbf{x} + a \leq 0\} \\
&= P\{\xi \geq x_1 > -a\} - P\{\xi < x_1 \leq -a\} \\
&= \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} (-1)^{\mathbb{1}\{\xi < -a\}} f(\mathbf{x}) \, dx_1 \, d\mathbf{x}_{2:d} \\
&= \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} (-1)^{\mathbb{1}\{\xi < -a\}} (f(-a, \mathbf{x}_{2:d}) + (f(\mathbf{x}) - f(-a, \mathbf{x}_{2:d}))) \, dx_1 \, d\mathbf{x}_{2:d}.
\end{aligned}$$

Regarding the second summand,

$$\begin{aligned}
&\left| \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} (-1)^{\mathbb{1}\{\xi < -a\}} (f(\mathbf{x}) - f(-a, \mathbf{x}_{2:d})) \, dx_1 \, d\mathbf{x}_{2:d} \right| \\
&= \left| \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} (-1)^{\mathbb{1}\{\xi < -a\}} \int_{-a}^{x_1} \nabla f(\eta, \mathbf{x}_{2:d})^\top \mathbf{e}_1 \, d\eta \, dx_1 \, d\mathbf{x}_{2:d} \right| \\
&\leq \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} \int_{-a}^{x_1} \|\nabla f(\eta, \mathbf{x}_{2:d})\| \, d\eta \, dx_1 \, d\mathbf{x}_{2:d} \\
&\leq \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} \int_{I(\xi)} \|\nabla f(\eta, \mathbf{x}_{2:d})\| \, d\eta \, dx_1 \, d\mathbf{x}_{2:d} \\
&= \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} |\xi + a| \|\nabla f(\eta, \mathbf{x}_{2:d})\| \, d\eta \, d\mathbf{x}_{2:d}.
\end{aligned}$$

3. Central Limit Theorems

Cauchy-Schwarz inequality yields

$$\begin{aligned}
 |\xi + a| &= \frac{n^{-1/2}}{|1 + n^{-1/2}\Delta_{\psi,1}|} \left| \Delta_{\psi}^{\top} \begin{pmatrix} a \\ -\mathbf{x}_{2:d} \\ -1 \end{pmatrix} \right| \\
 &\leq \frac{n^{-1/2}}{|1 + n^{-1/2}\Delta_{\psi,1}|} \|\Delta_{\psi}\| \sqrt{a^2 + 1} (1 + \|\mathbf{x}_{2:d}\|) \\
 &= M_n (1 + \|\mathbf{x}_{2:d}\|)
 \end{aligned}$$

with

$$M_n := \frac{n^{-1/2} \|\Delta_{\psi}\| \sqrt{a^2 + 1}}{|1 + n^{-1/2}\Delta_{\psi,1}|} = O_p(n^{-1/2})$$

and thus

$$\begin{aligned}
 &\left| \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} (-1)^{\mathbb{1}\{\xi < -a\}} (f(\mathbf{x}) - f(-a, \mathbf{x}_{2:d})) \, d\mathbf{x}_1 \, d\mathbf{x}_{2:d} \right| \\
 &\leq M_n \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} \|\nabla f(\eta, \mathbf{x}_{2:d})\| (1 + \|\mathbf{x}_{2:d}\|) \, d\eta \, d\mathbf{x}_{2:d} \\
 &\leq M_n \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} \sup_{\tilde{\eta} \in I(\xi)} \|\nabla f(\tilde{\eta}, \mathbf{x}_{2:d})\| (1 + \|\mathbf{x}_{2:d}\|) \, d\eta \, d\mathbf{x}_{2:d} \\
 &\leq M_n^2 \int_{\mathbb{R}^{d-1}} \sup_{\tilde{\eta} \in I(\xi)} \|\nabla f(\tilde{\eta}, \mathbf{x}_{2:d})\| (1 + \|\mathbf{x}_{2:d}\|)^2 \, d\mathbf{x}_{2:d} \\
 &\leq M_n^2 K = O_p(n^{-1})
 \end{aligned}$$

by Condition (3.3). Hence

$$\begin{aligned}
 &\sqrt{n} (P\{\mathbf{x} \in \mathbb{R}^d : \hat{\mathbf{b}}^{\top} \mathbf{x} + \hat{a} \leq 0\} - P\{\mathbf{x} \in \mathbb{R}^d : \mathbf{b}^{\top} \mathbf{x} + a \leq 0\}) \\
 &= \sqrt{n} \int_{\mathbb{R}^{d-1}} (a + \xi) f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} + O_p(n^{-1/2}) \\
 &= \int_{\mathbb{R}^{d-1}} \frac{\Delta_{\psi,1}a - \Delta_{\psi,d+1} - \sum_{i=2}^d \Delta_{\psi,i}x_i}{1 + n^{-\frac{1}{2}}\Delta_{\psi,1}} f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} + O_p(n^{-1/2}) \\
 &= \frac{\Delta_{\psi,1}a - \Delta_{\psi,d+1}}{1 + n^{-\frac{1}{2}}\Delta_{\psi,1}} \int_{\mathbb{R}^{d-1}} f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} \\
 &\quad - \frac{1}{1 + n^{-\frac{1}{2}}\Delta_{\psi,1}} \sum_{i=2}^d \Delta_{\psi,i} \int_{\mathbb{R}^{d-1}} x_i f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} + O_p(n^{-1/2})
 \end{aligned}$$

$$\begin{aligned}
&= -\frac{1}{1+n^{-\frac{1}{2}}\Delta_{\psi,1}} f_{X_1}(-a)\mathbb{E}((\mathbf{X}^\top, 1) \mid X_1 = -a)\Delta_{\psi} + O_p(n^{-1/2}) \\
&= -f_{X_1}(-a)\mathbb{E}((\mathbf{X}^\top, 1) \mid X_1 = -a)\Delta_{\psi} + O_p(n^{-1/2}).
\end{aligned}$$

Now for an arbitrary vector $\mathbf{b} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$, $\mathbf{b}^\top \mathbf{x} + a \leq 0$ if and only if $\mathbf{e}_1^\top \mathbf{y} + a/\|\mathbf{b}\| \leq 0$ with $\mathbf{y} = \mathbf{B}^\top \mathbf{x}$ and $\mathbf{B} = [\mathbf{b}/\|\mathbf{b}\|, \mathbf{b}_2, \dots, \mathbf{b}_d]$ such that $\mathbf{B}^\top \mathbf{B} = \mathbb{I}_d$. Further $\widehat{\mathbf{b}}^\top \mathbf{x} + \widehat{a} \leq 0$ if and only if $\widehat{\mathbf{b}}^\top \mathbf{y} + \widehat{a}' \leq 0$ with $\widehat{a}' = \widehat{a}/\|\mathbf{b}\|$ and $\widehat{\mathbf{b}}' = \mathbf{B}^\top \widehat{\mathbf{b}}/\|\mathbf{b}\|$. Then with $\mathbf{V} := \mathbf{B}^\top \mathbf{X}$,

$$\begin{aligned}
&P\{\mathbf{x} \in \mathbb{R}^d: \widehat{\mathbf{b}}^\top \mathbf{x} + \widehat{a} \leq 0\} - P\{\mathbf{x} \in \mathbb{R}^d: \mathbf{b}^\top \mathbf{x} + a \leq 0\} \\
&= -f_{V_1}(-a/\|\mathbf{b}\|)\mathbb{E}((\mathbf{V}^\top, 1) \mid V_1 = -a/\|\mathbf{b}\|) \begin{pmatrix} \widehat{\mathbf{b}}' - \mathbf{e}_1 \\ \widehat{a}' - a/\|\mathbf{b}\| \end{pmatrix} \\
&= -f_{\mathbf{b}^\top \mathbf{X}/\|\mathbf{b}\|}(-a/\|\mathbf{b}\|)\mathbb{E}(((\mathbf{B}^\top \mathbf{X})^\top, 1) \mid \mathbf{b}^\top \mathbf{X} = -a) \\
&\quad \cdot \frac{1}{\|\mathbf{b}\|} \begin{pmatrix} \mathbf{B}^\top & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{b}} - \mathbf{b} \\ \widehat{a} - a \end{pmatrix} \\
&= -f_{\mathbf{b}^\top \mathbf{X}}(-a)\mathbb{E}((\mathbf{X}^\top, 1) \mid \mathbf{b}^\top \mathbf{X} = -a) \begin{pmatrix} \widehat{\mathbf{b}} - \mathbf{b} \\ \widehat{a} - a \end{pmatrix}.
\end{aligned}$$

Next we show assertion (3.4). Without loss of generality let $\mathbf{b} = \mathbf{e}_1$. Then

$$\begin{aligned}
&P(\{\mathbf{x} \in \mathbb{R}^d: \widehat{\mathbf{b}}^\top \mathbf{x} + \widehat{a} \leq 0\} \triangle \{\mathbf{x} \in \mathbb{R}^d: \mathbf{b}^\top \mathbf{x} + a \leq 0\}) \\
&= P\{\xi \geq x_1 > -a\} + P\{\xi < x_1 \leq -a\} \\
&= \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} f(\mathbf{x}) \, d\mathbf{x}_1 \, d\mathbf{x}_{2:d} \\
&= \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} (f(-a, \mathbf{x}_{2:d}) + (f(\mathbf{x}) - f(-a, \mathbf{x}_{2:d}))) \, d\mathbf{x}_1 \, d\mathbf{x}_{2:d} \\
&= \int_{\mathbb{R}^{d-1}} |a + \xi| f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} + O_p(n^{-1}) \\
&= \frac{1}{\sqrt{n}} \int_{\mathbb{R}^{d-1}} \left| \frac{\Delta_{\psi,1}a - \Delta_{\psi,d+1} - \sum_{i=2}^d \Delta_{\psi,i}x_i}{1 + n^{-\frac{1}{2}}\Delta_{\psi,1}} \right| f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} \\
&\quad + O_p(n^{-1}).
\end{aligned}$$

3. Central Limit Theorems

In the penultimate step we used that

$$\begin{aligned}
& \left| \int_{\mathbb{R}^{d-1}} \int_{I(\xi)} (f(\mathbf{x}) - f(-a, \mathbf{x}_{2:d})) \, dx_1 \, d\mathbf{x}_{2:d} \right| \\
& \leq \int_{\mathbb{R}^{d-1}} \int_{I(\xi) - a}^{x_1} \|\nabla f(\eta, \mathbf{x}_{2:d})\| \, d\eta \, dx_1 \, d\mathbf{x}_{2:d} \\
& \leq M_n^2 K = O_p(n^{-1}).
\end{aligned}$$

Thus

$$\begin{aligned}
& \sqrt{n}P(\{\mathbf{x} \in \mathbb{R}^d: \hat{\mathbf{b}}^\top \mathbf{x} + \hat{a} \leq 0\} \triangle \{\mathbf{x} \in \mathbb{R}^d: \mathbf{b}^\top \mathbf{x} + a \leq 0\}) \\
& \leq \left| \frac{\Delta_{\psi,1}a - \Delta_{\psi,d+1}}{1 + n^{-\frac{1}{2}}\Delta_{\psi,1}} \right| \int_{\mathbb{R}^{d-1}} f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} \\
& \quad + \sum_{i=2}^d \left(\left| \frac{\Delta_{\psi,i}}{1 + n^{-\frac{1}{2}}\Delta_{\psi,1}} \right| \int_{\mathbb{R}^{d-1}} x_i f(-a, \mathbf{x}_{2:d}) \, d\mathbf{x}_{2:d} \right) + O_p(n^{-1/2}) \\
& = O_p(1).
\end{aligned}$$

□

3.1.2. Empirical Processes

For the proofs of the central limit theorems we need some results about empirical processes. The *empirical measure* \hat{P} of a sample of independent random variables $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ in a measurable space $(\mathcal{X}, \mathcal{B})$ with distribution P is the discrete random measure given by

$$\hat{P} := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{X}_i},$$

with $\delta_{\mathbf{x}}$ denoting the Dirac measure at \mathbf{x} . For a measurable set $D \subset \mathcal{X}$,

$$\hat{P}(D) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\mathbf{X}_i \in D\}.$$

For a probability measure P on \mathbb{R}^d we consider the empirical process

$$\mathbb{B}_{P,n} = \mathbb{B}_{P,n}(H)_{H \in \mathcal{H}}$$

with

$$\mathbb{B}_{P,n}(H) := \sqrt{n}(\hat{P} - P)(H),$$

where $\mathcal{H} := \{H(\beta, \gamma) : \beta \in \mathbb{R}^d \setminus \{\mathbf{0}\}, \gamma \in \mathbb{R}\}$. Instead of $\mathbb{B}_{P,n}(H(\beta, \gamma))$ we also write $\mathbb{B}_{P,n}(\beta, \gamma)$. For the proof of the following theorem we refer to van der Vaart and Wellner (1996). It is a consequence of the fact that \mathcal{H} is a Vapnik–Červonenkis class.

Theorem 3.4. *The empirical process $\mathbb{B}_{P,n}$ converges in $l_\infty(\mathcal{H})$ weakly to a centered Gaussian process B_P with covariances*

$$\mathbb{E}(B_P(H)B_P(H')) = P(H \cap H') - P(H)P(H')$$

and uniformly continuous sample paths with respect to

$$\rho_P(H, H') := P(H \triangle H').$$

Moreover, \mathcal{H} equipped with $\rho_P(H, H')$ is totally bounded and $\|\mathbb{B}_{P,n}\|_\infty = O_p(1)$, where $\|\cdot\|_\infty$ denotes the uniform norm.

Theorem 3.5. *Let $H(\beta, \gamma), H(\hat{\beta}, \hat{\gamma}) \in \mathcal{H}$ and P a probability measure on \mathbb{R}^d such that $P(H(\beta, \gamma) \triangle H(\hat{\beta}, \hat{\gamma})) \rightarrow_p 0$. Then*

$$\mathbb{B}_{P,n}(\hat{\beta}, \hat{\gamma}) - \mathbb{B}_{P,n}(\beta, \gamma) \rightarrow_p 0.$$

PROOF. For any fixed $\varepsilon > 0$, Theorem 3.4 implies that with asymptotic probability 1,

$$\begin{aligned} |\mathbb{B}_n(\hat{\beta}, \hat{\gamma}) - \mathbb{B}_n(\beta, \gamma)| &\leq \sup_{\substack{H(\tilde{\beta}, \tilde{\gamma}) \in \mathcal{H}: \\ \rho_P(H(\beta, \gamma), H(\tilde{\beta}, \tilde{\gamma})) \leq \varepsilon}} |\mathbb{B}_n(\tilde{\beta}, \tilde{\gamma}) - \mathbb{B}_n(\beta, \gamma)| \\ &\rightarrow_{\mathcal{L}} \sup_{\substack{H(\tilde{\beta}, \tilde{\gamma}) \in \mathcal{H}: \\ \rho_P(H(\beta, \gamma), H(\tilde{\beta}, \tilde{\gamma})) \leq \varepsilon}} |B_P(\tilde{\beta}, \tilde{\gamma}) - B_P(\beta, \gamma)| \\ &\rightarrow_p 0 \end{aligned}$$

as $\varepsilon \rightarrow 0$. □

3.2. Asymptotics of Estimators for Location and Scatter

The results of this section are not restricted to $L = 2$. For notational convenience we assume without loss of generality that $Y_1 = 1, Y_2 = 2, \dots, Y_L = L$.

Lemma 3.6. *Let $\theta \in \mathcal{Y}$ and $\hat{\mu}_\theta = N_\theta^{-1} \sum_{i \in \mathcal{G}_\theta} \mathbf{X}_i$ be the standard estimator for $\mu_\theta = \mathbb{E}(\mathbf{X}_i | Y_i = \theta)$. Suppose that $\mathbb{E}(\|\mathbf{X}\|^2) < \infty$. Then*

$$\Delta_{\mu_\theta} := \sqrt{n}(\hat{\mu}_\theta - \mu_\theta) = \sum_{i=1}^n \mathbf{Y}_{n,i}^{\mu_\theta} \rightarrow_{\mathcal{L}} \mathcal{N}_d(\mathbf{0}, w_\theta^{-1} \Sigma),$$

where $\mathbf{Y}_{n,i}^{\mu_\theta} := \sqrt{n} N_\theta^{-1} \mathbf{1}\{Y_i = \theta\} \tilde{\mathbf{X}}_i$.

3. Central Limit Theorems

Lemma 3.7. *Let $\widehat{\Sigma} := (n - L)^{-1} \sum_{\theta \in \mathcal{Y}} \sum_{i \in \mathcal{G}_\theta} (\mathbf{X}_i - \widehat{\boldsymbol{\mu}}_\theta)(\mathbf{X}_i - \widehat{\boldsymbol{\mu}}_\theta)^\top$ be the standard estimator for Σ and suppose that $\mathbb{E}(\|\mathbf{X}\|^4) < \infty$. Then*

$$\Delta_\Sigma := \sqrt{n}(\widehat{\Sigma} - \Sigma) = \sum_{i=1}^n \mathbf{Y}_{n,i}^\Sigma + o_p(1) \rightarrow_{\mathcal{L}} \mathcal{N}_{d \times d}(\mathbf{0}, \text{Var}(\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^\top)) \quad (3.6)$$

and

$$\begin{aligned} \Delta_{\Sigma^{-1}} &:= \sqrt{n}(\widehat{\Sigma}^{-1} - \Sigma^{-1}) = \sum_{i=1}^n \mathbf{Y}_{n,i}^{\Sigma^{-1}} + o_p(1) \\ &\rightarrow_{\mathcal{L}} \mathcal{N}_{d \times d}(\mathbf{0}, (\Sigma^{-1} \otimes \Sigma^{-1}) \text{Var}(\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^\top) (\Sigma^{-1} \otimes \Sigma^{-1})) \end{aligned} \quad (3.7)$$

with $\mathbf{Y}_{n,i}^\Sigma := n^{-1/2}(\tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^\top - \Sigma)$ and $\mathbf{Y}_{n,i}^{\Sigma^{-1}} := -\Sigma^{-1} \mathbf{Y}_{n,i}^\Sigma \Sigma^{-1}$.

Dümbgen et al. (2013) showed that the assumptions of the following lemma are satisfied for the M -estimators defined in Section 1.1.2 if $\mathcal{L}(\mathbf{X} \mid Y = \theta)$ is elliptically symmetric and $\mathbb{E}(\|\mathbf{X}\|^2) < \infty$.

Lemma 3.8. *Let*

$$\begin{aligned} \widehat{\Sigma} &= \Sigma + n^{-1} \sum_{\theta \in \mathcal{Y}} \sum_{i \in \mathcal{G}_\theta} \left(g_\theta(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_i\|) \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^\top - h_\theta(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_i\|) \Sigma \right) \\ &\quad + o_p(n^{-1/2}) \end{aligned}$$

for continuous bounded functions g_θ, h_θ such that $g_\theta(r)r^2$ is bounded for $r \geq 0$. Suppose that for $Y_i = \theta$

$$\mathbb{E}(g_\theta(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_i\|) \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^\top - h_\theta(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_i\|) \Sigma) = \mathbf{0}.$$

Then

$$\Delta_\Sigma = \sum_{i=1}^n \mathbf{Y}_{n,i}^\Sigma + o_p(1) \rightarrow_{\mathcal{L}} \mathcal{N}_{d \times d}(\mathbf{0}, V) \quad (3.8)$$

and

$$\Delta_{\Sigma^{-1}} = \sum_{i=1}^n \mathbf{Y}_{n,i}^{\Sigma^{-1}} + o_p(1) \rightarrow_{\mathcal{L}} \mathcal{N}_{d \times d}(\mathbf{0}, (\Sigma^{-1} \otimes \Sigma^{-1}) V (\Sigma^{-1} \otimes \Sigma^{-1})), \quad (3.9)$$

where

$$\begin{aligned} \mathbf{Y}_{n,i}^\Sigma &:= n^{-1/2} (g_{Y_i}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_i\|) \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^\top - h_{Y_i}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_i\|) \Sigma), \\ \mathbf{Y}_{n,i}^{\Sigma^{-1}} &:= -\Sigma^{-1} \mathbf{Y}_{n,i}^\Sigma \Sigma^{-1} \end{aligned}$$

and

$$\mathbf{V} := \sum_{\theta \in \mathcal{Y}} w_{\theta} \text{Var} \left(g_{\theta}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_{\theta}\|) \tilde{\mathbf{X}}_{\theta} \tilde{\mathbf{X}}_{\theta}^{\top} - h_{\theta}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_{\theta}\|) \Sigma \right).$$

PROOF OF LEMMA 3.6. We apply the central limit theorem (Theorem A.1) to the vectors $\mathbf{Y}_{n,i}^{\mu_{\theta}}$. The assumptions are fulfilled, since $\mathbb{E}(\mathbf{Y}_{n,i}^{\mu_{\theta}}) = \mathbf{0}$,

$$\sum_{i=1}^n \text{Var}(\mathbf{Y}_{n,i}^{\mu_{\theta}}) = N_{\theta} \text{Var}(\mathbf{Y}_{n,\theta}^{\mu_{\theta}}) = \frac{n}{N_{\theta}} \text{Var}(\mathbf{X}_1) = \frac{1}{\hat{w}_{\theta}} \Sigma \rightarrow \frac{1}{w_{\theta}} \Sigma$$

and by dominated convergence,

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}(\|\mathbf{Y}_{n,i}^{\mu_{\theta}}\|^2 \min(1, \|\mathbf{Y}_{n,i}^{\mu_{\theta}}\|)) &= \frac{n}{N_{\theta}} \mathbb{E}(\|\tilde{\mathbf{X}}_{\theta}\|^2 \min(1, \sqrt{n} N_{\theta}^{-1} \|\tilde{\mathbf{X}}_{\theta}\|)) \\ &\rightarrow \mathbf{0}. \end{aligned}$$

□

PROOF OF LEMMA 3.7. First note that

$$\begin{aligned} \hat{\Sigma} - \Sigma &= \frac{1}{n-L} \sum_{i=1}^n (\tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^{\top} - \Sigma) + \frac{L}{n-L} \Sigma \\ &\quad - \sum_{\theta \in \mathcal{Y}} \frac{N_{\theta}}{n-L} (\hat{\mu}_{\theta} - \mu_{\theta}) (\hat{\mu}_{\theta} - \mu_{\theta})^{\top} \\ &= n^{-1} \sum_{i=1}^n (\tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^{\top} - \Sigma) + O_p(n^{-1}). \end{aligned}$$

We apply the central limit theorem to $\sum_{i=1}^n \mathbf{Y}_{n,i}^{\Sigma} = \Delta_{\Sigma} + o_p(1)$. As to the assumptions, note that $\mathbb{E}(\mathbf{Y}_{n,i}^{\Sigma}) = \mathbf{0}$,

$$\sum_{i=1}^n \text{Var}(\mathbf{Y}_{n,i}^{\Sigma}) = n \text{Var}(\mathbf{Y}_{n,1}^{\Sigma}) = \text{Var}(\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^{\top})$$

and by dominated convergence,

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}(\|\mathbf{Y}_{n,i}^{\Sigma}\|^2 \min(1, \|\mathbf{Y}_{n,i}^{\Sigma}\|)) &= \mathbb{E}(\|\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^{\top} - \Sigma\|^2 \min(1, n^{-1/2} \|\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^{\top} - \Sigma\|)) \\ &\rightarrow \mathbf{0}. \end{aligned}$$

3. Central Limit Theorems

For an invertible matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ and $\mathbb{R}^{n \times n} \ni \mathbf{\Delta} \rightarrow \mathbf{0}$, it is well-known (e.g. Taylor and Lay, 1980) that

$$(\mathbf{B} + \mathbf{\Delta})^{-1} = \mathbf{B}^{-1} - \mathbf{B}^{-1} \mathbf{\Delta} \mathbf{B}^{-1} + O(\|\mathbf{\Delta}\|^2).$$

Consequently,

$$\hat{\Sigma}^{-1} = \Sigma^{-1} - n^{-1/2} \Sigma^{-1} \mathbf{\Delta}_{\Sigma} \Sigma^{-1} + O_p(n^{-1})$$

and

$$\mathbf{\Delta}_{\Sigma^{-1}} = -\Sigma^{-1} \mathbf{\Delta}_{\Sigma} \Sigma^{-1} + O_p(n^{-1/2}).$$

For arbitrary matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$ such that \mathbf{ABC} is well-defined, the relation

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}^{\top} \otimes \mathbf{A}) \text{vec}(\mathbf{B}) \quad (3.10)$$

is well-known and yields

$$\text{vec}(\Sigma^{-1} \mathbf{\Delta}_{\Sigma} \Sigma^{-1}) = (\Sigma^{-1} \otimes \Sigma^{-1}) \text{vec}(\mathbf{\Delta}_{\Sigma}).$$

Now with $\mathbf{Y}_{n,i}^{\Sigma^{-1}} := -\Sigma^{-1} \mathbf{Y}_{n,i}^{\Sigma} \Sigma^{-1}$, Claim (3.7) follows from (3.6). \square

PROOF OF LEMMA 3.8. By assumption, $\mathbb{E}(\mathbf{Y}_{n,i}^{\Sigma}) = \mathbf{0}$. Moreover, since $\sqrt{n} \|\mathbf{Y}_{n,i}^{\Sigma}\|$ is uniformly bounded,

$$\begin{aligned} \sum_{i=1}^n \text{Var}(\mathbf{Y}_{n,i}^{\Sigma}) &= \sum_{\theta \in \mathcal{Y}} N_{\theta} \text{Var}(\mathbf{Y}_{n,\theta}^{\Sigma}) \\ &= \sum_{\theta \in \mathcal{Y}} \frac{N_{\theta}}{n} \text{Var}(g_{\theta}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_{\theta}\|) \tilde{\mathbf{X}}_{\theta} \tilde{\mathbf{X}}_{\theta}^{\top} - h_{\theta}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_{\theta}\|) \Sigma) \\ &\rightarrow \sum_{\theta \in \mathcal{Y}} w_{\theta} \text{Var}(g_{\theta}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_{\theta}\|) \tilde{\mathbf{X}}_{\theta} \tilde{\mathbf{X}}_{\theta}^{\top} - h_{\theta}(\|\Sigma^{-1/2} \tilde{\mathbf{X}}_{\theta}\|) \Sigma) \end{aligned}$$

and for some constant $M < \infty$

$$\sum_{i=1}^n \mathbb{E}(\|\mathbf{Y}_{n,i}^{\Sigma}\|^2 \min(1, \|\mathbf{Y}_{n,i}^{\Sigma}\|)) \leq n \mathbb{E}(n^{-1} M^2 \cdot \min(1, n^{-1/2} M)) \rightarrow 0.$$

Now the central limit theorem yields (3.8). Claim (3.9) can be proved the same way as Claim (3.7) in Lemma 3.7. \square

3.3. A Central Limit Theorem for Missclassification Rates

In this section we examine the asymptotic properties of the standard linear classifier (1.3) for two classes, without assuming Gaussian distributions. More precisely, we describe the asymptotic distribution of the missclassification rates

$$R_\theta = \mathbb{P}(\hat{Y}(\mathbf{X}, \mathcal{D}) \neq Y \mid Y = \theta, \mathcal{D})$$

and the cross-validated estimators

$$\hat{R}_\theta = N_\theta^{-1} \#\{i \in \mathcal{G}_\theta : \hat{Y}(\mathbf{X}_i, \mathcal{D}_i) \neq Y_i\}$$

thereof.

Since we don't want to make assumptions on the convergence rate of \hat{w}_θ , we consider a reweighted version of the optimal classifier depending on \hat{w}_θ instead of w_θ , namely

$$\hat{Y}_n^*(\mathbf{x}) := \begin{cases} 1, & (\mathbf{x} - \boldsymbol{\mu}_{1,2})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \log(\hat{w}_2/\hat{w}_1) \leq 0 \\ 2, & (\mathbf{x} - \boldsymbol{\mu}_{1,2})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \log(\hat{w}_2/\hat{w}_1) > 0. \end{cases}$$

Note that we consider the class labels as fixed and therefore \hat{w}_θ is deterministic and converges to w_θ by assumption.

We define $\gamma := \|\boldsymbol{\beta}\|/2 + \log(w_1/w_2)/\|\boldsymbol{\beta}\|$, $\boldsymbol{\nu} := \mathbb{E}(\mathbf{Z} \mid \hat{Y}^*(\mathbf{X}_1) = 1)$, $\mathbf{u} := \|\boldsymbol{\beta}\|^{-1}\boldsymbol{\beta}$ and suppose that

$$\|\mathbb{E}(\mathbf{X}_1 \mid (\mathbf{X}_1 - \boldsymbol{\mu}_{1,2})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \log(\hat{w}_2/\hat{w}_1) = 0)\| < \infty, \quad (3.11)$$

where we condition on \mathbf{X}_1 being in the separating hyperplane of the optimal classifier.

The following two central limit theorems imply that, under certain conditions, the standard linear classifier for two classes is asymptotically optimal. Moreover, \hat{R}_1 is a root- n -consistent estimator of R_1 .

First, we consider the standard estimator of $\boldsymbol{\Sigma}$. In this case, the only assumptions we need to make about the distributions are a finite fourth moment and that (3.1) and (3.11) hold.

Theorem 3.9. *Let $L = 2$ and $\hat{\boldsymbol{\Sigma}}$ be the standard estimator. Suppose that $\mathbb{E}(\|\mathbf{X}\|^4) < \infty$, (3.1) and (3.11) hold. Then for the standard linear classifier,*

$$\sqrt{n} \begin{pmatrix} R_1 - \mathbb{P}(\hat{Y}_n^*(\mathbf{X}) \neq Y \mid Y = 1) \\ \hat{R}_1 - \mathbb{P}(\hat{Y}_n^*(\mathbf{X}) \neq Y \mid Y = 1) \end{pmatrix} \rightarrow_{\mathcal{L}} \mathcal{N}_2(\mathbf{0}, \boldsymbol{\Lambda})$$

3. Central Limit Theorems

and the components of the covariance matrix $\mathbf{\Lambda}$ are given by

$$\begin{aligned}\Lambda_{1,1} &= \|\boldsymbol{\beta}\|^{-2} f_{\mathbf{u}^\top \mathbf{Z}}(\gamma)^2 \left[\text{Var} \left(\boldsymbol{\beta}^\top \mathbf{Z} (\boldsymbol{\nu} - 2^{-1} \boldsymbol{\beta})^\top \mathbf{Z} \right) + w_1^{-1} \|\boldsymbol{\nu}\|^2 \right. \\ &\quad \left. + w_2^{-1} \|\boldsymbol{\nu} - \boldsymbol{\beta}\|^2 + 2\mathbb{E}((\boldsymbol{\beta}^\top \mathbf{Z})^2 (\boldsymbol{\nu} - 2^{-1} \boldsymbol{\beta})^\top \mathbf{Z}) \right], \\ \Lambda_{2,2} &= \Lambda_{1,1} + w_1^{-1} (\mathbb{P}(\widehat{Y}^*(\mathbf{X}_1) = 1) - \mathbb{P}(\widehat{Y}^*(\mathbf{X}_1) = 1)^2) \\ &\quad + 2\|\boldsymbol{\beta}\|^{-1} f_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \text{Cov} \left(\boldsymbol{\beta}^\top \mathbf{Z} (\boldsymbol{\nu} - 2^{-1} \boldsymbol{\beta})^\top \mathbf{Z}, \mathbb{1}\{\widehat{Y}^*(\mathbf{X}_1) = 1\} \right) \\ &\quad + 2w_1^{-1} \|\boldsymbol{\beta}\|^{-1} f_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \boldsymbol{\nu}^\top \mathbb{E}(\mathbf{Z} \mathbb{1}\{\widehat{Y}^*(\mathbf{X}_1) = 1\})\end{aligned}$$

and

$$\begin{aligned}\Lambda_{1,2} &= \Lambda_{2,1} \\ &= \Lambda_{1,1} + \|\boldsymbol{\beta}\|^{-1} f_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \text{Cov} \left(\boldsymbol{\beta}^\top \mathbf{Z} (\boldsymbol{\nu} - 2^{-1} \boldsymbol{\beta})^\top \mathbf{Z}, \mathbb{1}\{\widehat{Y}^*(\mathbf{X}_1) = 1\} \right) \\ &\quad + w_1^{-1} \|\boldsymbol{\beta}\|^{-1} f_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \boldsymbol{\nu}^\top \mathbb{E}(\mathbf{Z} \mathbb{1}\{\widehat{Y}^*(\mathbf{X}_1) = 1\}).\end{aligned}$$

Instead of the standard estimator for the covariance matrix $\boldsymbol{\Sigma}$, one could use the more robust M -estimators defined in Section 1.1.2. Dümbgen et al. (2013) showed that these estimators satisfy the assumptions of Lemma 3.8 if $\mathcal{L}(\mathbf{Z})$ is spherically symmetric. In this case Condition (3.11) is not necessary and Condition (3.1) can be relaxed to

$$K(\mathbf{U}) := \int_{\mathbb{R}^{d-1}} \sup_{\mathbf{a} \in \mathbf{U}^\perp} \|\nabla f_{\mathbf{Z}}(\mathbf{a} + \mathbf{U}\mathbf{z})\| \, d\mathbf{z} < \infty. \quad (3.12)$$

Theorem 3.10. *Let $L = 2$ and $\widehat{\boldsymbol{\Sigma}}$ be an estimator satisfying the assumptions of Lemma 3.8. Suppose that $\mathcal{L}(\mathbf{Z})$ is spherically symmetric satisfying (3.12) and $\mathbb{E}(\|\mathbf{Z}\|^2) < \infty$. Then for the standard linear classifier,*

$$\sqrt{n} \begin{pmatrix} R_1 - \mathbb{P}(\widehat{Y}_n^*(\mathbf{X}) \neq Y \mid Y = 1) \\ \widehat{R}_1 - \mathbb{P}(\widehat{Y}_n^*(\mathbf{X}) \neq Y \mid Y = 1) \end{pmatrix} \rightarrow_{\mathcal{L}} \mathcal{N}_2(\mathbf{0}, \mathbf{\Lambda})$$

and the components of the covariance matrix $\mathbf{\Lambda}$ are given by

$$\begin{aligned}\Lambda_{1,1} &= f_{Z_1}(\gamma)^2 \left[w_1(2^{-1} \|\boldsymbol{\beta}\| - \gamma)^2 \text{Var} \left(g_1(\|\mathbf{Z}\|) Z_1^2 - h_1(\|\mathbf{Z}\|) \right) \right. \\ &\quad + w_2(2^{-1} \|\boldsymbol{\beta}\| - \gamma)^2 \text{Var} \left(g_2(\|\mathbf{Z}\|) Z_1^2 - h_2(\|\mathbf{Z}\|) \right) \\ &\quad + w_1^{-1} (\gamma / \|\boldsymbol{\beta}\|)^2 + w_2^{-1} (\gamma / \|\boldsymbol{\beta}\| - 1)^2 \\ &\quad + \gamma(2\gamma / \|\boldsymbol{\beta}\| - 1) \mathbb{E} \left(g_1(\|\mathbf{Z}\|) Z_1^3 - h_1(\|\mathbf{Z}\|) Z_1 \right) \\ &\quad \left. + (\|\boldsymbol{\beta}\| - \gamma)(2\gamma / \|\boldsymbol{\beta}\| - 1) \mathbb{E} \left(g_2(\|\mathbf{Z}\|) Z_1^3 - h_2(\|\mathbf{Z}\|) Z_1 \right) \right],\end{aligned}$$

$$\begin{aligned}
 \Lambda_{2,2} &= \Lambda_{1,1} + w_1^{-1} (\mathbb{P}(Z_1 \leq \gamma) - \mathbb{P}(Z_1 \leq \gamma)^2) \\
 &\quad + f_{Z_1}(\gamma)(2\gamma - \|\beta\|)\mathbb{E}\left((g_1(\|\mathbf{Z}\|)Z_1^2 - h_1(\|\mathbf{Z}\|))\mathbb{1}\{Z_1 \leq \gamma\}\right) \\
 &\quad + 2w_1^{-1}f_{Z_1}(\gamma)(\gamma/\|\beta\|)\mathbb{E}(Z_1\mathbb{1}\{Z_1 \leq \gamma\})
 \end{aligned}$$

and

$$\begin{aligned}
 \Lambda_{1,2} &= \Lambda_{2,1} \\
 &= \Lambda_{1,1} + f_{Z_1}(\gamma)(\gamma - \|\beta\|/2)\mathbb{E}\left((g_1(\|\mathbf{Z}\|)Z_1^2 - h_1(\|\mathbf{Z}\|))\mathbb{1}\{Z_1 \leq \gamma\}\right) \\
 &\quad + w_1^{-1}f_{Z_1}(\gamma)(\gamma/\|\beta\|)\mathbb{E}(Z_1\mathbb{1}\{Z_1 \leq \gamma\}).
 \end{aligned}$$

PROOF OF THEOREM 3.9. Note that

$$\begin{aligned}
 &\mathbb{P}(\widehat{Y}_n^*(\mathbf{X}) = 1 \mid Y = 1) \\
 &= \mathbb{P}((\mathbf{X} - \boldsymbol{\mu}_{1,2})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \log(\widehat{w}_2/\widehat{w}_1) \leq 0 \mid Y = 1) \\
 &= P_1\{\mathbf{x} \in \mathbb{R}^d : \mathbf{b}^\top \mathbf{x} + a_n \leq 0\},
 \end{aligned}$$

where $a_n := (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{1,2} + \log(\widehat{w}_2/\widehat{w}_1)$, $\mathbf{b} := \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$ and

$$\begin{aligned}
 &\mathbb{P}(\widehat{Y}(\mathbf{X}, \mathcal{D}) = 1 \mid \mathcal{D}, Y = 1) \\
 &= \mathbb{P}((\mathbf{X} - \widehat{\boldsymbol{\mu}}_{1,2})^\top \widehat{\boldsymbol{\Sigma}}^{-1}(\widehat{\boldsymbol{\mu}}_2 - \widehat{\boldsymbol{\mu}}_1) + \log(\widehat{w}_2/\widehat{w}_1) \leq 0 \mid \mathcal{D}, Y = 1) \\
 &= P_1\{\mathbf{x} \in \mathbb{R}^d : \widehat{\mathbf{b}}^\top \mathbf{x} + \widehat{a} \leq 0\},
 \end{aligned}$$

where $\widehat{a} := (\widehat{\boldsymbol{\mu}}_1 - \widehat{\boldsymbol{\mu}}_2)^\top \widehat{\boldsymbol{\Sigma}}^{-1} \widehat{\boldsymbol{\mu}}_{1,2} + \log(\widehat{w}_2/\widehat{w}_1)$ and $\widehat{\mathbf{b}} := \widehat{\boldsymbol{\Sigma}}^{-1}(\widehat{\boldsymbol{\mu}}_2 - \widehat{\boldsymbol{\mu}}_1)$ with estimators $\widehat{\boldsymbol{\mu}}_\theta(\mathcal{D})$ and $\widehat{\boldsymbol{\Sigma}}^{-1}(\mathcal{D})$. Elementary calculations reveal that $\widehat{\boldsymbol{\mu}}_\theta(\mathcal{D}_i(\mathbf{X})) = \widehat{\boldsymbol{\mu}}_\theta(\mathcal{D}) + \text{O}_p(n^{-1})$ and $\widehat{\boldsymbol{\Sigma}}^{-1}(\mathcal{D}_i(\mathbf{X})) = \widehat{\boldsymbol{\Sigma}}^{-1}(\mathcal{D}) + \text{O}_p(n^{-1})$. Thus

$$\begin{aligned}
 \widehat{R}_1 &= \frac{\#\{i \in \mathcal{G}_1 : \widehat{\mathbf{b}}(\mathcal{D}_i)^\top \mathbf{X}_i + \widehat{a}(\mathcal{D}_i) > 0\}}{N_1} \\
 &= \frac{\#\{i \in \mathcal{G}_1 : \widehat{\mathbf{b}}(\mathcal{D})^\top \mathbf{X}_i + \widehat{a}(\mathcal{D}) + \text{O}_p(n^{-1}) > 0\}}{N_1} \\
 &= \widehat{P}_1\{\mathbf{x} \in \mathbb{R}^d : \widehat{\mathbf{b}}(\mathcal{D})^\top \mathbf{x} + \widehat{a}' > 0\}
 \end{aligned}$$

with $\widehat{a}' = \widehat{a}(\mathcal{D}) + \text{O}_p(n^{-1})$.

Define $H_n := \{\mathbf{x} \in \mathbb{R}^d : \mathbf{b}^\top \mathbf{x} + a_n \leq 0\}$, $\widehat{H} := \{\mathbf{x} \in \mathbb{R}^d : \widehat{\mathbf{b}}^\top \mathbf{x} + \widehat{a} \leq 0\}$ and

3. Central Limit Theorems

$\hat{H}' := \{\mathbf{x} \in \mathbb{R}^d : \hat{\mathbf{b}}^\top \mathbf{x} + \hat{a}' \leq 0\}$. Then

$$\begin{aligned}\boldsymbol{\eta}_n &:= \sqrt{n} \begin{pmatrix} P_1(\hat{H}) - P_1(H_n) \\ \hat{P}_1(\hat{H}') - P_1(H_n) \end{pmatrix} \\ &= -\sqrt{n} \begin{pmatrix} R_1 - \mathbb{P}(\hat{Y}_n^*(\mathbf{X}) \neq Y \mid Y = 1) \\ \hat{R}_1 - \mathbb{P}(\hat{Y}_n^*(\mathbf{X}) \neq Y \mid Y = 1) \end{pmatrix}.\end{aligned}$$

Next note that

$$\begin{aligned}\hat{\mathbf{b}} &= \hat{\boldsymbol{\Sigma}}^{-1}(\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) \\ &= (\boldsymbol{\Sigma}^{-1} + n^{-1/2} \boldsymbol{\Delta}_{\boldsymbol{\Sigma}^{-1}})((\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + n^{-1/2}(\boldsymbol{\Delta}_{\boldsymbol{\mu}_2} - \boldsymbol{\Delta}_{\boldsymbol{\mu}_1})) \\ &= \mathbf{b} + n^{-1/2}(\boldsymbol{\Delta}_{\boldsymbol{\Sigma}^{-1}}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \boldsymbol{\Sigma}^{-1}(\boldsymbol{\Delta}_{\boldsymbol{\mu}_2} - \boldsymbol{\Delta}_{\boldsymbol{\mu}_1})) + \text{O}_p(n^{-1}).\end{aligned}$$

Hence

$$\begin{aligned}\boldsymbol{\Delta}_{\mathbf{b}} &:= \sqrt{n}(\hat{\mathbf{b}} - \mathbf{b}) \\ &= \boldsymbol{\Delta}_{\boldsymbol{\Sigma}^{-1}}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \boldsymbol{\Sigma}^{-1}(\boldsymbol{\Delta}_{\boldsymbol{\mu}_2} - \boldsymbol{\Delta}_{\boldsymbol{\mu}_1}) + \text{O}_p(n^{-1/2}) \\ &= \sum_{i=1}^n \mathbf{Y}_{n,i}^{\mathbf{b}} + \text{o}_p(1),\end{aligned}$$

where $\mathbf{Y}_{n,i}^{\mathbf{b}} := \mathbf{Y}_{n,i}^{\boldsymbol{\Sigma}^{-1}}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \boldsymbol{\Sigma}^{-1}(\mathbf{Y}_{n,i}^{\boldsymbol{\mu}_2} - \mathbf{Y}_{n,i}^{\boldsymbol{\mu}_1})$ with $\mathbf{Y}_{n,i}^{\boldsymbol{\Sigma}^{-1}}$ and $\mathbf{Y}_{n,i}^{\boldsymbol{\mu}_\theta}$ as in Lemma 3.6 and 3.7. Moreover,

$$\begin{aligned}\hat{a} &= (\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2)^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_{1,2} + \log(\hat{w}_2/\hat{w}_1) \\ &= a_n + n^{-1/2}((\boldsymbol{\Delta}_{\boldsymbol{\mu}_1} - \boldsymbol{\Delta}_{\boldsymbol{\mu}_2})^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{1,2} + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Delta}_{\boldsymbol{\Sigma}^{-1}} \boldsymbol{\mu}_{1,2} \\ &\quad + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Delta}_{\boldsymbol{\mu}_{1,2}}) + \text{O}_p(n^{-1})\end{aligned}$$

and

$$\begin{aligned}\Delta_a &:= \sqrt{n}(\hat{a} - a_n) \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Delta}_{\boldsymbol{\Sigma}^{-1}} \boldsymbol{\mu}_{1,2} + \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Delta}_{\boldsymbol{\mu}_1} - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Delta}_{\boldsymbol{\mu}_2} + \text{O}_p(n^{-1/2}) \\ &= \sum_{i=1}^n Y_{n,i}^a + \text{o}_p(1),\end{aligned}$$

where $Y_{n,i}^a := (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{Y}_{n,i}^{\boldsymbol{\Sigma}^{-1}} \boldsymbol{\mu}_{1,2} + \boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}^{-1} \mathbf{Y}_{n,i}^{\boldsymbol{\mu}_1} - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}^{-1} \mathbf{Y}_{n,i}^{\boldsymbol{\mu}_2}$.

Define $\boldsymbol{\psi}_n^\top := (\mathbf{b}^\top, a_n)^\top$ and

$$\boldsymbol{\Delta}_{\boldsymbol{\psi}} := \sqrt{n}(\hat{\boldsymbol{\psi}} - \boldsymbol{\psi}_n) = \begin{pmatrix} \boldsymbol{\Delta}_{\mathbf{b}} \\ \Delta_a \end{pmatrix} = \sum_{i=1}^n \mathbf{Y}_{n,i}^{\boldsymbol{\psi}} + \text{o}_p(1)$$

with

$$\mathbf{Y}_{n,i}^\psi := \begin{pmatrix} \mathbf{Y}_{n,i}^b \\ \mathbf{Y}_{n,i}^a \end{pmatrix}.$$

By Lemma 3.1,

$$\begin{aligned} \sqrt{n}(P_1(\hat{H}) - P_1(H_n)) &= \sqrt{n}(P_1(\hat{H}') - P_1(H_n)) + o_p(1) \\ &= \mathbf{c}^\top \boldsymbol{\Delta}_\psi + o_p(1) \\ &= \sum_{i=1}^n \mathbf{c}^\top \mathbf{Y}_{n,i}^\psi + o_p(1), \end{aligned}$$

where

$$\begin{aligned} \mathbf{c} = \mathbf{c}_n &= -f_{\mathbf{b}^\top \mathbf{X}_1}(-a_n) \mathbb{E} \left(\begin{pmatrix} \mathbf{X}_1 \\ 1 \end{pmatrix} \middle| \mathbf{b}^\top \mathbf{X}_1 = -a_n \right)^\top \\ &= -f_{\mathbf{b}^\top \mathbf{X}_1}(-a_n) \begin{pmatrix} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\nu}_n + \boldsymbol{\mu}_1 \\ 1 \end{pmatrix} \end{aligned}$$

and $\boldsymbol{\nu}_n := \mathbb{E}(\mathbf{Z} \mid \mathbf{b}^\top \mathbf{X}_1 = -a_n)$.

Regarding the second component of $\boldsymbol{\eta}_n$, note that

$$\hat{P}_1(\hat{H}') - P_1(H_n) = (P_1(\hat{H}') - P_1(H_n)) + (\hat{P}_1 - P_1)(H_n) + R_n$$

with

$$\begin{aligned} R_n &:= (\hat{P}_1 - P_1)(\hat{H}') - (\hat{P}_1 - P_1)(H_n) \\ &= n^{-1/2} (\mathbb{B}_{P_1,n}(\hat{\mathbf{b}}, \hat{a}') - \mathbb{B}_{P_1,n}(\mathbf{b}, a_n)) \end{aligned}$$

and the empirical process $\mathbb{B}_{P_1,n}$ defined in Section 3.1.2. Lemma 3.1 entails that $P_1(\hat{H}' \triangle H_n) \rightarrow_p 0$ and we deduce from Theorem 3.5 that $R_n = o_p(n^{-1/2})$. Next we define

$$Y_{n,i}^p := \frac{\mathbb{1}\{i \in \mathcal{G}_1\}}{\sqrt{n\hat{w}_1}} (\mathbb{1}\{\mathbf{X}_i \in H_n\} - P_1(H_n))$$

such that

$$\begin{aligned} \sum_{i=1}^n Y_{n,i}^p &= \frac{\sqrt{n}}{N_1} \sum_{i \in \mathcal{G}_1} (\mathbb{1}\{\mathbf{X}_i \in H_n\} - P_1(H_n)) \\ &= \sqrt{n}(\hat{P}_1 - P_1)(H_n) \end{aligned}$$

3. Central Limit Theorems

and

$$\begin{aligned}\eta_n &= \sqrt{n} \left(\begin{array}{c} P_1(\hat{H}) - P_1(H_n) \\ (P_1(\hat{H}') - P_1(H_n)) + (\hat{P}_1 - P_1)(H_n) \end{array} \right) + o_p(1) \\ &= \sum_{i=1}^n \mathbf{Y}_{n,i}^\eta + o_p(1)\end{aligned}$$

with

$$\mathbf{Y}_{n,i}^\eta := \begin{pmatrix} \mathbf{c}^\top \mathbf{Y}_{n,i}^\psi \\ \mathbf{c}^\top \mathbf{Y}_{n,i}^\psi + Y_{n,i}^p \end{pmatrix}.$$

Before we can apply the central limit theorem to $\sum_{i=1}^n \mathbf{Y}_{n,i}^\eta$ we have to compute $\text{Var}(\mathbf{Y}_{n,i}^\eta)$, the sum $\mathbf{\Lambda}^n := \sum_{i=1}^n \text{Var}(\mathbf{Y}_{n,i}^\eta)$ and its limit $\mathbf{\Lambda} = \lim_{n \rightarrow \infty} \mathbf{\Lambda}^n$. To this end note that

$$\begin{aligned}\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi &= \mathbf{c}_{1:d}^\top \mathbf{Y}_{n,i}^b + c_{d+1} Y_{n,i}^a \\ &= \mathbf{c}_{1:d}^\top \mathbf{Y}_{n,i}^{\Sigma^{-1}} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \mathbf{c}_{1:d}^\top \Sigma^{-1} (\mathbf{Y}_{n,i}^{\mu_2} - \mathbf{Y}_{n,i}^{\mu_1}) \\ &\quad + c_{d+1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{Y}_{n,i}^{\Sigma^{-1}} \boldsymbol{\mu}_{1,2} + c_{d+1} \boldsymbol{\mu}_1^\top \Sigma^{-1} \mathbf{Y}_{n,i}^{\mu_1} \\ &\quad - c_{d+1} \boldsymbol{\mu}_2^\top \Sigma^{-1} \mathbf{Y}_{n,i}^{\mu_2} \\ &= (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{1,2})^\top \mathbf{Y}_{n,i}^{\Sigma^{-1}} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \\ &\quad + (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{Y_i})^\top \Sigma^{-1} (\mathbf{Y}_{n,i}^{\mu_2} - \mathbf{Y}_{n,i}^{\mu_1}) \\ &= ((\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \otimes (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{1,2}))^\top \text{vec}(\mathbf{Y}_{n,i}^{\Sigma^{-1}}) \\ &\quad + \frac{\sqrt{n}}{N_{Y_i}} (-1)^{\mathbb{1}\{Y_i=1\}} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{Y_i})^\top \Sigma^{-1} \tilde{\mathbf{X}}_i \\ &= -\mathbf{v}_n^\top \text{vec}(\mathbf{Y}_{n,i}^\Sigma) + \frac{\sqrt{n}}{N_{Y_i}} (-1)^{\mathbb{1}\{Y_i=1\}} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{Y_i})^\top \Sigma^{-1} \tilde{\mathbf{X}}_i,\end{aligned}$$

where

$$\begin{aligned}\mathbf{v}_n &:= (\Sigma^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)) \otimes (\Sigma^{-1}(\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{1,2})) \in \mathbb{R}^{d^2} \\ &= -\mathbf{f}_{\mathbf{b}^\top \mathbf{X}_1}(-a_n) \mathbf{b} \otimes (\Sigma^{-1/2}(\boldsymbol{\nu}_n - 2^{-1} \boldsymbol{\beta})).\end{aligned}$$

Here and for the following computations we utilize several times the relations

(1.22), (1.23) and (3.10). The upper left component of $\mathbf{\Lambda}$ is given by

$$\begin{aligned}
 \Lambda_{1,1}^n &= \sum_{i=1}^n \text{Var}(\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi) \\
 &= \sum_{i=1}^n \left(\text{Var}(\mathbf{v}_n^\top \text{vec}(\mathbf{Y}_{n,i}^\Sigma)) \right. \\
 &\quad + \frac{n}{N_{Y_i}^2} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{Y_i})^\top \Sigma^{-1} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{Y_i}) \\
 &\quad \left. + \frac{2\sqrt{n}(-1)^{\mathbb{1}\{Y_i=2\}}}{N_{Y_i}} \mathbf{v}_n^\top \text{Cov}(\text{vec}(\mathbf{Y}_{n,i}^\Sigma), \tilde{\mathbf{X}}_i) \Sigma^{-1} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{Y_i}) \right) \\
 &= \text{Var}(\mathbf{v}_n^\top \text{vec}(\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^\top)) + f_{\mathbf{b}^\top \mathbf{X}_1}(-a_n)^2 (\hat{w}_1^{-1} \|\boldsymbol{\nu}_n\|^2 + \hat{w}_2^{-1} \|\boldsymbol{\nu}_n - \boldsymbol{\beta}\|^2) \\
 &\quad + 2c_{d+1} \mathbf{v}_n^\top \text{Cov}(\text{vec}(\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^\top), \tilde{\mathbf{X}}_1) \Sigma^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \\
 &\rightarrow \|\boldsymbol{\beta}\|^{-2} f_{\mathbf{u}^\top \mathbf{Z}}(\gamma)^2 \left[\text{Var}(\boldsymbol{\beta}^\top \mathbf{Z}(\boldsymbol{\nu} - 2^{-1} \boldsymbol{\beta})^\top \mathbf{Z}) + w_1^{-1} \|\boldsymbol{\nu}\|^2 \right. \\
 &\quad \left. + w_2^{-1} \|\boldsymbol{\nu} - \boldsymbol{\beta}\|^2 + 2\mathbb{E}((\boldsymbol{\beta}^\top \mathbf{Z})^2 (\boldsymbol{\nu} - 2^{-1} \boldsymbol{\beta})^\top \mathbf{Z}) \right],
 \end{aligned}$$

where $\mathbf{u} := \|\boldsymbol{\beta}\|^{-1} \boldsymbol{\beta}$. We used that $f_{\mathbf{b}^\top \mathbf{X}_1}(-a) = \|\boldsymbol{\beta}\|^{-1} f_{\mathbf{u}^\top \mathbf{Z}}(\gamma)$ and $\boldsymbol{\nu}_n \rightarrow \boldsymbol{\nu}$. The latter assertion can be derived from (3.1).

To compute $\text{Var}(\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi + Y_{n,i}^p)$ we need

$$\text{Var}(Y_{n,i}^p) = \frac{\mathbb{1}\{i \in \mathcal{G}_1\}}{n\hat{w}_1^2} (P_1(H_n) - P_1(H_n)^2)$$

and

$$\begin{aligned}
 &\text{Cov}(\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi, Y_{n,i}^p) \\
 &= \text{Cov}(-\mathbf{v}_n^\top \text{vec}(\mathbf{Y}_{n,i}^\Sigma), Y_{n,i}^p) \\
 &\quad + \text{Cov}\left(\frac{\sqrt{n}}{N_{Y_i}} (-1)^{\mathbb{1}\{Y_i=1\}} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_{Y_i})^\top \Sigma^{-1} \tilde{\mathbf{X}}_i, Y_{n,i}^p\right) \\
 &= -\frac{\mathbb{1}\{i \in \mathcal{G}_1\}}{N_1} \mathbf{v}_n^\top \text{Cov}(\text{vec}(\tilde{\mathbf{X}}_1 \tilde{\mathbf{X}}_1^\top), \mathbb{1}\{\mathbf{X}_1 \in H_n\}) \\
 &\quad - \frac{\mathbb{1}\{i \in \mathcal{G}_1\}}{N_1 \hat{w}_1} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_1)^\top \Sigma^{-1} \text{Cov}(\tilde{\mathbf{X}}_1, \mathbb{1}\{\mathbf{X}_1 \in H_n\}) \\
 &= \frac{\mathbb{1}\{i \in \mathcal{G}_1\}}{N_1} f_{\mathbf{b}^\top \mathbf{X}_1}(-a_n) \text{Cov}(\boldsymbol{\beta}^\top \mathbf{Z}(\boldsymbol{\nu}_n - 2^{-1} \boldsymbol{\beta})^\top \mathbf{Z}, \mathbb{1}\{\mathbf{X}_1 \in H_n\}) \\
 &\quad + \frac{\mathbb{1}\{i \in \mathcal{G}_1\}}{N_1 \hat{w}_1} f_{\mathbf{b}^\top \mathbf{X}_1}(-a_n) \boldsymbol{\nu}_n^\top \mathbb{E}(\mathbf{Z} \mathbb{1}\{\mathbf{X}_1 \in H_n\}).
 \end{aligned}$$

3. Central Limit Theorems

Employing (3.1) again, one can show that $P_1(H_n) \rightarrow \mathbb{P}(\hat{Y}^*(\mathbf{X}_1) = 1)$. Thus $\text{Cov}(\mathbf{Z}, \mathbb{1}\{\mathbf{X}_1 \in H\}) \rightarrow \text{Cov}(\mathbf{Z}, \mathbb{1}\{\hat{Y}^*(\mathbf{X}_1) = 1\})$ and $\text{Cov}(\beta^\top \mathbf{Z}(\boldsymbol{\nu}_n - 2^{-1}\beta)^\top \mathbf{Z}, \mathbb{1}\{\mathbf{X}_1 \in H\}) \rightarrow \text{Cov}(\beta^\top \mathbf{Z}(\boldsymbol{\nu} - 2^{-1}\beta)^\top \mathbf{Z}, \mathbb{1}\{\hat{Y}^*(\mathbf{X}_1) = 1\})$ by dominated convergence. Hence the lower left component of $\mathbf{\Lambda}$ is given by

$$\begin{aligned} \Lambda_{2,2}^n &= \sum_{i=1}^n \text{Var}(\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi + Y_{n,i}^p) \\ &\rightarrow \Lambda_{1,1} + w_1^{-1} (\mathbb{P}(\hat{Y}^*(\mathbf{X}_1) = 1) - \mathbb{P}(\hat{Y}^*(\mathbf{X}_1) = 1)^2) \\ &\quad + 2\|\beta\|^{-1} \mathbf{f}_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \text{Cov}(\beta^\top \mathbf{Z}(\boldsymbol{\nu} - 2^{-1}\beta)^\top \mathbf{Z}, \mathbb{1}\{\hat{Y}^*(\mathbf{X}_1) = 1\}) \\ &\quad + 2w_1^{-1} \|\beta\|^{-1} \mathbf{f}_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \boldsymbol{\nu}^\top \mathbb{E}(\mathbf{Z} \mathbb{1}\{\hat{Y}^*(\mathbf{X}_1) = 1\}) \end{aligned}$$

and the other components are equal to

$$\begin{aligned} \Lambda_{1,2}^n &= \Lambda_{2,1}^n \\ &= \sum_{i=1}^n (\text{Var}(\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi) + \text{Cov}(\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi, Y_{n,i}^p)) \\ &\rightarrow \Lambda_{1,1} + \|\beta\|^{-1} \mathbf{f}_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \text{Cov}(\beta^\top \mathbf{Z}(\boldsymbol{\nu} - 2^{-1}\beta)^\top \mathbf{Z}, \mathbb{1}\{\hat{Y}^*(\mathbf{X}_1) = 1\}) \\ &\quad + w_1^{-1} \|\beta\|^{-1} \mathbf{f}_{\mathbf{u}^\top \mathbf{Z}}(\gamma) \boldsymbol{\nu}^\top \mathbb{E}(\mathbf{Z} \mathbb{1}\{\hat{Y}^*(\mathbf{X}_1) = 1\}). \end{aligned}$$

Next we show that Lindeberg's condition is satisfied. Note that

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}(\|\mathbf{Y}_{n,i}^\eta\|^2 \min(1, \|\mathbf{Y}_{n,i}^\eta\|)) &= \sum_{\lambda \in \{1,2\}} \sum_{i \in \mathcal{G}_\lambda} \mathbb{E}(\|\mathbf{Y}_{n,i}^\eta\|^2 \min(1, \|\mathbf{Y}_{n,i}^\eta\|)) \\ &= \sum_{\lambda \in \{1,2\}} N_\lambda \mathbb{E}(\|\mathbf{Y}_{n,\lambda}^\eta\|^2 \min(1, \|\mathbf{Y}_{n,\lambda}^\eta\|)) \end{aligned}$$

and for $\lambda \in \{1, 2\}$, $\|\mathbf{Y}_{n,\lambda}^\eta\| \leq 2\|\mathbf{c}\| \|\mathbf{Y}_{n,\lambda}^\psi\| + |Y_{n,\lambda}^p| \leq 2\|\mathbf{c}\| \|\mathbf{Y}_{n,\lambda}^\psi\| + n^{-1/2} \hat{w}_1^{-1}$. For some constants c' and c''

$$\|\mathbf{Y}_{n,\lambda}^\psi\| \leq \|\mathbf{Y}_{n,\lambda}^b\| + |Y_{n,\lambda}^a| \leq c' \|\mathbf{Y}_{n,\lambda}^\Sigma\| + c'' \frac{\sqrt{n}}{N_\lambda} \|\tilde{\mathbf{X}}_\lambda\|.$$

Therefore $\|\mathbf{Y}_{n,\lambda}^\Sigma\| \rightarrow_{\text{a.s.}} 0$ implies $\|\mathbf{Y}_{n,\lambda}^\psi\| \rightarrow_{\text{a.s.}} 0$ and $\|\mathbf{Y}_{n,\lambda}^\eta\| \rightarrow_{\text{a.s.}} 0$. Since

$$\mathbb{E}\|\tilde{\mathbf{X}}_\lambda\|^2 < \infty, \quad \mathbb{E}\sqrt{n}\|\mathbf{Y}_{n,\lambda}^\Sigma\| \|\tilde{\mathbf{X}}_\lambda\| < \infty \quad \text{and} \quad \mathbb{E}n\|\mathbf{Y}_{n,\lambda}^\Sigma\|^2 < \infty$$

by assumption, the dominated convergence theorem implies that

$$N_\lambda \mathbb{E}(\|\mathbf{Y}_{n,\lambda}^\eta\|^2 \min(1, \|\mathbf{Y}_{n,\lambda}^\eta\|)) \rightarrow 0$$

3.3. A Central Limit Theorem for Missclassification Rates

for any $\lambda \in \mathcal{Y}$ and thus Lindeberg's condition

$$\sum_{i=1}^n \mathbb{E} \left(\|\mathbf{Y}_{n,i}^\eta\|^2 \min(1, \|\mathbf{Y}_{n,i}^\eta\|) \right) \rightarrow 0$$

is satisfied. Since $\mathbb{E}(\mathbf{Y}_{n,i}^\eta) = \mathbf{0}$ for all $i \leq n$ and the Gaussian distribution is symmetric, the assertion follows from the central limit theorem (Theorem A.1). \square

PROOF OF THEOREM 3.10. The proof of this theorem is similar to the proof of Theorem 3.9. But the covariance matrix of the limit distribution is slightly different. The elliptical symmetry implies that

$$\begin{aligned} \mathbb{P}(\widehat{Y}^*(\mathbf{X}) = 1 \mid Y = 1) &= \mathbb{P}((\mathbf{X} - \boldsymbol{\mu}_{1,2})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \log(\widehat{w}_2/\widehat{w}_1) \leq 0 \mid Y = 1) \\ &= P_1\{\mathbf{x} \in \mathbb{R}^d: \mathbf{b}^\top \mathbf{x} + a_n \leq 0\} \\ &= P_1\{\boldsymbol{\beta}^\top \boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_1) + \alpha_n \leq 0\} \\ &= P_0\{\boldsymbol{\beta}^\top \mathbf{x} + \alpha_n \leq 0\} \\ &= P_0\{(\boldsymbol{\beta}/\|\boldsymbol{\beta}\|)^\top \mathbf{x} + \alpha_n/\|\boldsymbol{\beta}\| \leq 0\} \\ &= P_0\{x_1 + \alpha_n/\|\boldsymbol{\beta}\| \leq 0\}, \end{aligned}$$

where $\alpha_n = a_n + \mathbf{b}^\top \boldsymbol{\mu}_1$. Analogously we get

$$\mathbb{P}(\widehat{Y}(\mathbf{X}, \mathcal{D}) = 1 \mid \mathcal{D}, Y = 1) = P_0\{x_1 + \widehat{\alpha}/\|\widehat{\boldsymbol{\beta}}\| \leq 0\}$$

with $\widehat{\boldsymbol{\beta}} = \boldsymbol{\Sigma}^{1/2}\widehat{\mathbf{b}}$ and $\widehat{\alpha} = \widehat{a} + \widehat{\mathbf{b}}^\top \boldsymbol{\mu}_1$.

Employing Lemma 3.1 with $\widehat{\mathbf{b}} = \mathbf{b} = \mathbf{e}_1$, $a = \alpha_n/\|\boldsymbol{\beta}\|$ and $\widehat{a} = \widehat{\alpha}/\|\widehat{\boldsymbol{\beta}}\| = \alpha_n/\|\boldsymbol{\beta}\| + O_p(n^{-1/2})$ yields

$$\begin{aligned} \sqrt{n}(P_1(\widehat{H}) - P_1(H_n)) &= \mathbf{c}^\top \boldsymbol{\Delta}_\psi + o_p(1) \\ &= -f_{Z_1}(\gamma_n)\sqrt{n}\left(\frac{\widehat{\alpha}}{\|\widehat{\boldsymbol{\beta}}\|} - \frac{\alpha_n}{\|\boldsymbol{\beta}\|}\right) + o_p(1), \end{aligned}$$

where $\gamma_n := \|\boldsymbol{\beta}\|/2 + \log(\widehat{w}_1/\widehat{w}_2)/\|\boldsymbol{\beta}\| = -\alpha_n/\|\boldsymbol{\beta}\|$. Note that we use Lemma 3.1 with $\widehat{\mathbf{b}} = \mathbf{b} = \mathbf{e}_1$. Therefore Condition (3.2) is not necessary and Condition (3.3) can be relaxed to (3.12).

The first order Taylor expansion of $\|\widehat{\boldsymbol{\beta}}\|$ is given by

$$\|\widehat{\boldsymbol{\beta}}\| = \|\boldsymbol{\beta}\| + \frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \boldsymbol{\beta}}{\|\boldsymbol{\beta}\|} + o_p(n^{-1/2}).$$

3. Central Limit Theorems

Together with $\hat{\alpha} = \alpha_n + O_p(n^{-1/2})$ this entails that

$$\begin{aligned} \frac{\hat{\alpha}}{\|\hat{\beta}\|} - \frac{\hat{\alpha}}{\|\beta\|} &= \frac{-(\|\hat{\beta}\| - \|\beta\|)\alpha_n}{(\|\hat{\beta}\| - \|\beta\|)\|\beta\| + \|\beta\|^2} + o_p(n^{-1/2}) \\ &= \frac{-(\|\hat{\beta}\| - \|\beta\|)\alpha_n}{\|\beta\|^2} + o_p(n^{-1/2}) \\ &= -(\hat{\beta} - \beta)^\top \beta \frac{\alpha_n}{\|\beta\|^3} + o_p(n^{-1/2}) \end{aligned}$$

and thus

$$\begin{aligned} \sqrt{n} \left(\frac{\hat{\alpha}}{\|\hat{\beta}\|} - \frac{\alpha_n}{\|\beta\|} \right) &= -\sqrt{n}(\hat{\beta} - \beta)^\top \beta \frac{\alpha_n}{\|\beta\|^3} + \sqrt{n} \frac{\hat{\alpha} - \alpha_n}{\|\beta\|} + o_p(1) \\ &= -\frac{\alpha_n}{\|\beta\|^3} (\mu_2 - \mu_1)^\top \Delta_b + \frac{1}{\|\beta\|} \Delta_a + \frac{1}{\|\beta\|} \mu_1^\top \Delta_b + o_p(1) \\ &= \left(\left(\frac{\alpha_n}{\|\beta\|^3} + \frac{1}{\|\beta\|} \right) \mu_1 - \frac{\alpha_n}{\|\beta\|^3} \mu_2 \right)^\top \Delta_b + \frac{1}{\|\beta\|} \Delta_a + o_p(1). \end{aligned}$$

Therefore in the elliptic symmetric case \mathbf{c} is given by

$$\begin{aligned} \mathbf{c} &= -\frac{f_{Z_1}(\gamma_n)}{\|\beta\|} \begin{pmatrix} (\frac{\alpha_n}{\|\beta\|^2} + 1) \mu_1 - \frac{\alpha_n}{\|\beta\|^2} \mu_2 \\ 1 \end{pmatrix} \\ &= -\frac{f_{Z_1}(\gamma_n)}{\|\beta\|} \begin{pmatrix} (1 - \frac{\gamma_n}{\|\beta\|}) \mu_1 + \frac{\gamma_n}{\|\beta\|} \mu_2 \\ 1 \end{pmatrix} \end{aligned}$$

and

$$\mathbf{v}_n = \frac{f_{Z_1}(\gamma_n)}{\|\beta\|} \left(\frac{1}{2} - \frac{\gamma_n}{\|\beta\|} \right) (\mathbf{b} \otimes \mathbf{b}).$$

Hence the upper left component of $\mathbf{\Lambda}$ is given by

$$\begin{aligned} \Lambda_{1,1} &= f_{Z_1}(\gamma)^2 \left[w_1(2^{-1}\|\beta\| - \gamma)^2 \text{Var} \left(g_1(\|\mathbf{Z}\|) Z_1^2 - h_1(\|\mathbf{Z}\|) \right) \right. \\ &\quad + w_2(2^{-1}\|\beta\| - \gamma)^2 \text{Var} \left(g_2(\|\mathbf{Z}\|) Z_1^2 - h_2(\|\mathbf{Z}\|) \right) \\ &\quad + w_1^{-1}(\gamma/\|\beta\|)^2 + w_2^{-1}(\gamma/\|\beta\| - 1)^2 \\ &\quad + \gamma(2\gamma/\|\beta\| - 1) \mathbb{E} \left(g_1(\|\mathbf{Z}\|) Z_1^3 - h_1(\|\mathbf{Z}\|) Z_1 \right) \\ &\quad \left. + (\|\beta\| - \gamma)(2\gamma/\|\beta\| - 1) \mathbb{E} \left(g_2(\|\mathbf{Z}\|) Z_1^3 - h_2(\|\mathbf{Z}\|) Z_1 \right) \right]. \end{aligned}$$

We utilized that $\|\beta\|^{-1}\beta^\top \mathbf{Z}$ has the same distribution as Z_1 . Note that $P_1(H_n) = \mathbb{P}(Z_1 \leq \gamma_n)$, $\mathbf{1}\{\mathbf{X}_1 \in H_n\} = \mathbf{1}\{\|\beta\|^{-1}\beta^\top \mathbf{Z} \leq \gamma_n\}$ and

$$\begin{aligned}
 & \text{Cov}(\mathbf{c}^\top \mathbf{Y}_{n,i}^\psi, Y_{n,i}^p) \\
 &= -\frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1} \mathbf{v}_n^\top \text{Cov}(\sqrt{n} \text{vec}(\mathbf{Y}_{n,1}^\Sigma), \mathbf{1}\{\mathbf{X}_1 \in H_n\}) \\
 &\quad - \frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1 \hat{w}_1} (\mathbf{c}_{1:d} - c_{d+1} \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1} \text{Cov}(\tilde{\mathbf{X}}_1, \mathbf{1}\{\mathbf{X}_1 \in H_n\}) \\
 &= \frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1} f_{Z_1}(\gamma_n) (\gamma_n - \|\beta\|/2) \mathbb{E}\left((g_1(\|\mathbf{Z}\|) Z_1^2 - h_1(\|\mathbf{Z}\|)) \mathbf{1}\{\mathbf{X}_1 \in H_n\}\right) \\
 &\quad + \frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1 \hat{w}_1} \frac{\gamma_n f_{Z_1}(\gamma_n)}{\|\beta\|} \mathbb{E}(Z_1 \mathbf{1}\{\mathbf{X}_1 \in H_n\}) \\
 &= \frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1} f_{Z_1}(\gamma_n) (\gamma_n - \|\beta\|/2) \mathbb{E}\left((g_1(\|\mathbf{Z}\|) Z_1^2 - h_1(\|\mathbf{Z}\|)) \mathbf{1}\{Z_1 \leq \gamma_n\}\right) \\
 &\quad + \frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1 \hat{w}_1} \frac{\gamma_n f_{Z_1}(\gamma_n)}{\|\beta\|} \mathbb{E}(Z_1 \mathbf{1}\{Z_1 \leq \gamma_n\}).
 \end{aligned}$$

Therefore the remaining components are given by

$$\begin{aligned}
 \Lambda_{2,2}^n &= \Lambda_{1,1}^n + \hat{w}_1^{-1} (P_1(H_n) - P_1(H_n)^2) \\
 &\quad + \frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1} f_{Z_1}(\gamma_n) (2\gamma_n - \|\beta\|) \\
 &\quad \cdot \mathbb{E}\left((g_1(\|\mathbf{Z}\|) Z_1^2 - h_1(\|\mathbf{Z}\|)) \mathbf{1}\{Z_1 \leq \gamma_n\}\right) \\
 &\quad + 2 \frac{\mathbf{1}\{i \in \mathcal{G}_1\}}{N_1 \hat{w}_1} \frac{\gamma_n f_{Z_1}(\gamma_n)}{\|\beta\|} \mathbb{E}(Z_1 \mathbf{1}\{Z_1 \leq \gamma_n\}). \\
 &\rightarrow \Lambda_{1,1} + w_1^{-1} (\mathbb{P}(Z_1 \leq \gamma) - \mathbb{P}(Z_1 \leq \gamma)^2) \\
 &\quad + f_{Z_1}(\gamma) (2\gamma - \|\beta\|) \mathbb{E}\left((g_1(\|\mathbf{Z}\|) Z_1^2 - h_1(\|\mathbf{Z}\|)) \mathbf{1}\{Z_1 \leq \gamma\}\right) \\
 &\quad + 2w_1^{-1} f_{Z_1}(\gamma) (\gamma/\|\beta\|) \mathbb{E}(Z_1 \mathbf{1}\{Z_1 \leq \gamma\})
 \end{aligned}$$

and

$$\begin{aligned}
 \Lambda_{1,2} &= \Lambda_{2,1} \\
 &= \Lambda_{1,1} + f_{Z_1}(\gamma) (\gamma - \|\beta\|/2) \mathbb{E}\left((g_1(\|\mathbf{Z}\|) Z_1^2 - h_1(\|\mathbf{Z}\|)) \mathbf{1}\{Z_1 \leq \gamma\}\right) \\
 &\quad + w_1^{-1} f_{Z_1}(\gamma) (\gamma/\|\beta\|) \mathbb{E}(Z_1 \mathbf{1}\{Z_1 \leq \gamma\}).
 \end{aligned}$$

We used the dominated convergence theorem twice.

As to Lindeberg's condition, note that $\|\mathbf{Y}_{n,i}^\Sigma\| \leq n^{-1/2} M$ for some constant $M < \infty$. Thus the proof is similar to the one used for Theorem 3.9. \square

3.4. A Central Limit Theorem for Inclusion Probabilities

In this section we consider p-values based on the plug-in statistic for the standard model with two classes. The corresponding conditional inclusion probabilities

$$\mathcal{I}_\alpha(b, \theta \mid \mathcal{D}) = \mathbb{P}(\theta \in \hat{\mathcal{Y}}_\alpha(\mathbf{X}, \mathcal{D}) \mid Y = b, \mathcal{D})$$

are of interest to judge the separability of the two classes. However, these theoretic quantities are typically unknown. Therefore we use cross-validation to estimate them. Namely, we compute the empirical conditional inclusion probabilities

$$\hat{\mathcal{I}}_\alpha(b, \theta) = N_b^{-1} \# \{i \in \mathcal{G}_b : \theta \in \hat{\mathcal{Y}}_\alpha(\mathbf{X}_i, \mathcal{D}_i)\}$$

based on cross-validated p-values.

Dümbgen et al. (2008) showed that $\hat{\mathcal{I}}_\alpha(b, \theta)$ are consistent estimators of $\mathcal{I}_\alpha(b, \theta \mid \mathcal{D})$, (see also Section 1.6). More precisely,

$$\left. \begin{array}{l} \mathcal{I}_\alpha(b, \theta \mid \mathcal{D}) \\ \hat{\mathcal{I}}_\alpha(b, \theta) \end{array} \right\} \rightarrow_p \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = b),$$

and the limit equals $1 - \alpha$ in case of $b = \theta$.

We now take a closer look at the inclusion probabilities and describe the asymptotic distribution of

$$\sqrt{n} \begin{pmatrix} \mathcal{I}_\alpha(b, \theta \mid \mathcal{D}) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = b) \\ \hat{\mathcal{I}}_\alpha(b, \theta) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = b) \end{pmatrix}$$

assuming only elliptically symmetric instead of Gaussian distributions.

Let $E := \{\mathbf{x} \in \mathbb{R}^d : \pi_\theta^*(\mathbf{x}) > \alpha\}$ and q denote the $(1 - \alpha)$ -quantile of $\mathcal{L}(\mathbf{v}^\top \mathbf{Z})$ for some unit vector $\mathbf{v} \in \mathbb{R}^d$. The spherical symmetry of $\mathcal{L}(\mathbf{Z})$ implies that q does not depend on \mathbf{v} .

Theorem 3.11. *Suppose that $\mathcal{L}(\mathbf{Z})$ is elliptically symmetric satisfying (3.1) and $\mathbb{E}(\|\mathbf{Z}\|^2) < \infty$. Either let $\hat{\Sigma}$ be the standard estimator and $\mathbb{E}(\|\mathbf{Z}\|^4) < \infty$, or let $\hat{\Sigma}$ be an estimator satisfying the assumptions of Lemma 3.8, e.g. the M-estimators defined in Section 1.1.2. Then for the plug-in rule for the standard model with $L = 2$ classes and $\theta \in \{1, 2\}$,*

$$\begin{aligned} \sqrt{n} \left(\mathcal{I}_\alpha(\theta, \theta \mid \mathcal{D}) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = \theta) \right) &= -\sqrt{n}(\hat{P}_\theta - P_\theta)(E) + o_p(1) \\ &\rightarrow_{\mathcal{L}} \mathcal{N}(0, w_\theta^{-1} \alpha(1 - \alpha)) \end{aligned}$$

and

$$\sqrt{n} \left(\hat{\mathcal{I}}_\alpha(\theta, \theta) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = \theta) \right) = o_p(1). \quad (3.13)$$

For $\theta \neq \lambda \in \{1, 2\}$,

$$\begin{aligned} & \sqrt{n} \begin{pmatrix} \mathcal{I}_\alpha(\lambda, \theta \mid \mathcal{D}) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = \lambda) \\ \widehat{\mathcal{I}}_\alpha(\lambda, \theta) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = \lambda) \end{pmatrix} \\ &= \sqrt{n} \begin{pmatrix} c(\widehat{P}_\theta - P_\theta)(E) \\ c(\widehat{P}_\theta - P_\theta)(E) + (\widehat{P}_\lambda - P_\lambda)(E) \end{pmatrix} + o_p(1) \\ &\rightarrow_{\mathcal{L}} \mathcal{N}_2\left(\mathbf{0}, \begin{pmatrix} \sigma_1^2 & \sigma_1^2 \\ \sigma_1^2 & \sigma_1^2 + \sigma_2^2 \end{pmatrix}\right), \end{aligned}$$

where $\sigma_1^2 = w_\theta^{-1} c^2 \alpha (1 - \alpha)$ and $\sigma_2^2 = w_\lambda^{-1} P_\lambda(E) (1 - P_\lambda(E))$ with $c = -f_{Z_1}(q - \|\beta\|) / f_{Z_1}(q)$ and $P_\lambda(E) = F_{Z_1}(q - \|\beta\|)$.

Corollary 3.12. *Suppose that the assumptions of Theorem 3.11 are satisfied. Then for $b, \theta \in \{1, 2\}$*

$$\begin{aligned} \sqrt{n}(\widehat{\mathcal{I}}_\alpha(b, \theta) - \mathcal{I}_\alpha(b, \theta \mid \mathcal{D})) &= \sqrt{n}(\widehat{P}_b - P_b)(E) + o_p(1) \\ &\rightarrow_{\mathcal{L}} \mathcal{N}(0, w_b^{-1} P_b(E) (1 - P_b(E))). \end{aligned}$$

It is remarkable that the term of order $O_p(n^{-1/2})$ in $\widehat{\mathcal{I}}_\alpha(\theta, \theta) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = \theta)$ vanishes. This means that the cross-validated estimator $\widehat{\mathcal{I}}_\alpha(\theta, \theta)$ converges faster to the inclusion probability of the optimal p-value than the theoretic conditional inclusion probability $\mathcal{I}_\alpha(\theta, \theta \mid \mathcal{D})$ does.

The term of order $O_p(n^{-1/2})$ in $\widehat{\mathcal{I}}_\alpha(\lambda, \theta) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = \theta)$ does not vanish for $\lambda \neq \theta$. But it may be written as a sum of two independent summands. The first summand $\mathcal{I}_\alpha(\lambda, \theta \mid \mathcal{D})$ depends only on the training data from class λ and the second summand depends only on the training data from class θ .

The Corollary implies that the empirical conditional inclusion probabilities $\widehat{\mathcal{I}}_\alpha(b, \theta)$ are root- n -consistent estimators for the conditional inclusion probabilities $\mathcal{I}_\alpha(b, \theta \mid \mathcal{D})$. Moreover, it enables us to construct confidence intervals for $\mathcal{I}_\alpha(b, \theta \mid \mathcal{D})$. An asymptotic $(1 - \alpha)$ -confidence interval is given by

$$\left[\widehat{\mathcal{I}}_\alpha(b, \theta) \pm \sqrt{n^{-1} w_b^{-1} P_b(E) (1 - P_b(E)) z_{1-\alpha/2}} \right],$$

where $z_{1-\alpha/2}$ denotes the $(1 - \alpha/2)$ -quantile of the standard Gaussian distribution. In practice, the prior probability w_b may be unknown and can be replaced by the deterministic quantity $\widehat{w}_b = N_b/n$. For $b = \theta$, we know that $P_\theta(E) = \alpha$ and can compute the confidence interval

$$\left[\widehat{\mathcal{I}}_\alpha(\theta, \theta) \pm \sqrt{N_\theta^{-1} \alpha (1 - \alpha) z_{1-\alpha/2}} \right]$$

3. Central Limit Theorems

for $\mathcal{I}_\alpha(\theta, \theta \mid \mathcal{D})$. However, for $\lambda \neq \theta$, $P_\lambda(E)$ is typically unknown in practice, which prevents us from constructing confidence intervals for $\mathcal{I}_\alpha(\lambda, \theta \mid \mathcal{D})$.

For the proof of Theorem 3.11 we need the following lemma.

Lemma 3.13. *Suppose that the assumptions of Theorem 3.11 are satisfied. Then for $b, \theta \in \{1, 2\}$*

$$P_b(\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\} \triangle \{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) > \alpha\}) = O_p(n^{-1/2}) \quad (3.14)$$

and

$$\begin{aligned} & P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) > \alpha\} \\ &= \frac{f_{Z_1}(q - \|\Sigma^{-1/2}(\boldsymbol{\mu}_b - \boldsymbol{\mu}_\theta)\|)}{f_{Z_1}(q)}(\hat{P}_\theta - P_\theta)\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) \leq \alpha\} + o_p(n^{-1/2}). \end{aligned} \quad (3.15)$$

For $b = \theta$ this reduces to

$$\begin{aligned} & P_\theta\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\} - P_\theta\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) > \alpha\} \\ &= (\hat{P}_\theta - P_\theta)\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) \leq \alpha\} + o_p(n^{-1/2}). \end{aligned}$$

PROOF OF LEMMA 3.13. For $L = 2$ classes and $\theta \neq \lambda \in \{1, 2\}$, we consider

$$T_\theta^*(\mathbf{x}) := (\mathbf{x} - \boldsymbol{\mu}_{\lambda, \theta})^\top \Sigma^{-1}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta),$$

which is a strictly monotonic transformation of $T_\theta^*(\mathbf{x})$ defined in Example 1.2 and therefore leads to the same p-values. The empirical version of $T_\theta^*(\mathbf{x})$ based on training data \mathcal{D} is given by

$$T_\theta(\mathbf{x}, \mathcal{D}) = (\mathbf{x} - \hat{\boldsymbol{\mu}}_{\lambda, \theta})^\top \hat{\Sigma}^{-1}(\hat{\boldsymbol{\mu}}_\lambda - \hat{\boldsymbol{\mu}}_\theta),$$

with estimators $\hat{\boldsymbol{\mu}}_b(\mathcal{D})$ and $\hat{\Sigma}^{-1}(\mathcal{D})$. Elementary calculations reveal that $\hat{\boldsymbol{\mu}}_b(\mathcal{D}_i(\mathbf{X})) = \hat{\boldsymbol{\mu}}_b(\mathcal{D}) + O_p(n^{-1})$ and $\hat{\Sigma}^{-1}(\mathcal{D}_i(\mathbf{X})) = \hat{\Sigma}^{-1}(\mathcal{D}) + O_p(n^{-1})$. Thus $T_\theta(\mathbf{X}_i, \mathcal{D}_i(\mathbf{X})) = T_\theta(\mathbf{X}_i, \mathcal{D}) + O_p(n^{-1})$ and

$$\begin{aligned} \pi_\theta(\mathbf{X}, \mathcal{D}) &:= \frac{\#\{i \in \mathcal{G}_\theta: T_\theta(\mathbf{X}_i, \mathcal{D}_i(\mathbf{X})) \geq T_\theta(\mathbf{X}, \mathcal{D})\} + 1}{N_\theta + 1} \\ &= \hat{P}_\theta\{\mathbf{z} \in \mathbb{R}^d: T_\theta(\mathbf{z}, \mathcal{D}) + O_p(n^{-1}) \geq T_\theta(\mathbf{X}, \mathcal{D})\} + O_p(n^{-1}). \end{aligned}$$

Consequently,

$$\begin{aligned} & P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) > \alpha\} \\ &= P_b\left\{\mathbf{x} \in \mathbb{R}^d: \hat{P}_\theta\{\mathbf{z} \in \mathbb{R}^d: T_\theta(\mathbf{z}, \mathcal{D}) + O_p(n^{-1}) \geq T_\theta(\mathbf{x}, \mathcal{D})\} \right. \\ &\quad \left. + O_p(n^{-1}) > \alpha\right\} - P_b\left\{\mathbf{x} \in \mathbb{R}^d: P_\theta\{\mathbf{z} \in \mathbb{R}^d: T_\theta^*(\mathbf{z}) \geq T_\theta^*(\mathbf{x})\} > \alpha\right\} \\ &= P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{P}_\theta(\hat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(A(\mathbf{x})) > \alpha\} \end{aligned}$$

where

$$\begin{aligned}\widehat{A}(\mathbf{x}) &:= \{\mathbf{z} \in \mathbb{R}^d: T_\theta(\mathbf{z}, \mathcal{D}) + O_p(n^{-1}) \geq T_\theta(\mathbf{x}, \mathcal{D})\}, \\ A(\mathbf{x}) &:= \{\mathbf{z} \in \mathbb{R}^d: T_\theta^*(\mathbf{z}) \geq T_\theta^*(\mathbf{x})\}.\end{aligned}$$

Now we split the term in two summands

$$\begin{aligned}& P_b\{\mathbf{x} \in \mathbb{R}^d: \widehat{P}_\theta(\widehat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(A(\mathbf{x})) > \alpha\} \\ &= \left(P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(\widehat{A}(\mathbf{x})) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(A(\mathbf{x})) > \alpha\} \right) \\ &+ \left(P_b\{\mathbf{x} \in \mathbb{R}^d: \widehat{P}_\theta(\widehat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} \right. \\ &\quad \left. - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(\widehat{A}(\mathbf{x})) > \alpha\} \right).\end{aligned}\tag{3.16}$$

Regarding the first summand note that

$$\begin{aligned}A(\mathbf{x}) &= \{\mathbf{z} \in \mathbb{R}^d: T_\theta^*(\mathbf{z}) \geq T_\theta^*(\mathbf{x})\} \\ &= \{\mathbf{z} \in \mathbb{R}^d: \mathbf{z}^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta) \geq \mathbf{x}^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta)\} \\ &= \{\mathbf{z} \in \mathbb{R}^d: \boldsymbol{\beta}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{z} \geq \boldsymbol{\beta}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{x}\} \\ &= \{\mathbf{z} \in \mathbb{R}^d: \mathbf{u}^\top \boldsymbol{\Sigma}^{-1/2}(\mathbf{z} - \boldsymbol{\mu}_\theta) \geq \mathbf{u}^\top \boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)\},\end{aligned}$$

where $\boldsymbol{\beta} = \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta)$ and $\mathbf{u} := \boldsymbol{\beta}/\|\boldsymbol{\beta}\|$. Note that $P_\theta(A(\mathbf{x})) > \alpha$ if and only if $\mathbf{u}^\top \boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) < q$. Consequently,

$$P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(A(\mathbf{x})) > \alpha\} = P_b\{\mathbf{x} \in \mathbb{R}^d: \mathbf{u}^\top \boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) < q\}.\tag{3.17}$$

Lemma 3.6, 3.7 and 3.8 yield

$$\begin{aligned}T_\theta(\mathbf{x}, \mathcal{D}) &= (\mathbf{x} - \widehat{\boldsymbol{\mu}}_{\lambda, \theta})^\top \widehat{\boldsymbol{\Sigma}}^{-1}(\widehat{\boldsymbol{\mu}}_\lambda - \widehat{\boldsymbol{\mu}}_\theta) \\ &= (\mathbf{x} - \boldsymbol{\mu}_{\lambda, \theta})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta) - \frac{\boldsymbol{\Delta}_{\boldsymbol{\mu}_{\lambda, \theta}}^\top}{\sqrt{n}} \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta) \\ &\quad + (\mathbf{x} - \boldsymbol{\mu}_{\lambda, \theta})^\top \frac{\boldsymbol{\Delta}_{\boldsymbol{\Sigma}^{-1}}}{\sqrt{n}}(\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta) + (\mathbf{x} - \boldsymbol{\mu}_{\lambda, \theta})^\top \boldsymbol{\Sigma}^{-1} \frac{\boldsymbol{\Delta}_{\boldsymbol{\mu}_\lambda} - \boldsymbol{\Delta}_{\boldsymbol{\mu}_\theta}}{\sqrt{n}} \\ &\quad + O_p(n^{-1})\end{aligned}$$

3. Central Limit Theorems

and thus

$$\begin{aligned}
\widehat{A}(\mathbf{x}) &= \{ \mathbf{z} \in \mathbb{R}^d : T_\theta(\mathbf{z}, \mathcal{D}) + O_p(n^{-1}) \geq T_\theta(\mathbf{x}, \mathcal{D}) \} \\
&= \left\{ \mathbf{z} \in \mathbb{R}^d : T_\theta^*(\mathbf{z}) - T_\theta^*(\mathbf{x}) + (\mathbf{z} - \mathbf{x})^\top \Sigma^{-1/2} \left(\Sigma^{1/2} \frac{\Delta \Sigma^{-1}}{\sqrt{n}} (\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta) \right. \right. \\
&\quad \left. \left. + \Sigma^{-1/2} \frac{\Delta \boldsymbol{\mu}_\lambda - \Delta \boldsymbol{\mu}_\theta}{\sqrt{n}} + O_p(n^{-1}) \right) \geq 0 \right\} \\
&= \left\{ \mathbf{z} \in \mathbb{R}^d : (\mathbf{z} - \mathbf{x})^\top \Sigma^{-1/2} \left(\boldsymbol{\beta} + \Sigma^{1/2} \frac{\Delta \Sigma^{-1}}{\sqrt{n}} (\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta) \right. \right. \\
&\quad \left. \left. + \Sigma^{-1/2} \frac{\Delta \boldsymbol{\mu}_\lambda - \Delta \boldsymbol{\mu}_\theta}{\sqrt{n}} + O_p(n^{-1}) \right) \geq 0 \right\} \\
&= \{ \mathbf{z} \in \mathbb{R}^d : \widehat{\boldsymbol{\beta}}^\top \Sigma^{-1/2} \mathbf{z} \geq \widehat{\boldsymbol{\beta}}^\top \Sigma^{-1/2} \mathbf{x} \} \\
&= \{ \mathbf{z} \in \mathbb{R}^d : \widehat{\mathbf{u}}^\top \Sigma^{-1/2} (\mathbf{z} - \boldsymbol{\mu}_\theta) \geq \widehat{\mathbf{u}}^\top \Sigma^{-1/2} (\mathbf{x} - \boldsymbol{\mu}_\theta) \},
\end{aligned}$$

where $\widehat{\boldsymbol{\beta}} := \boldsymbol{\beta} + \Sigma^{1/2} \frac{\Delta \Sigma^{-1}}{\sqrt{n}} (\boldsymbol{\mu}_\lambda - \boldsymbol{\mu}_\theta) + \Sigma^{-1/2} \frac{\Delta \boldsymbol{\mu}_\lambda - \Delta \boldsymbol{\mu}_\theta}{\sqrt{n}} + O_p(n^{-1})$ and $\widehat{\mathbf{u}} := \widehat{\boldsymbol{\beta}} / \|\widehat{\boldsymbol{\beta}}\|$.

We deduce from the spherically symmetric distribution of $\Sigma^{-1/2}(\mathbf{X}_\theta - \boldsymbol{\mu}_\theta)$ that $\widehat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{X}_\theta - \boldsymbol{\mu}_\theta)$ conditional on the training data has the same distribution as $\mathbf{u}^\top \Sigma^{-1/2}(\mathbf{X}_\theta - \boldsymbol{\mu}_\theta)$. Therefore $P_\theta\{\mathbf{z} \in \mathbb{R}^d : T_\theta(\mathbf{z}, \mathcal{D}) + O_p(n^{-1}) \geq T_\theta(\mathbf{x}, \mathcal{D})\} > \alpha$ if and only if $\widehat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) < q$.

This entails that

$$P_b\{\mathbf{x} \in \mathbb{R}^d : P_\theta(\widehat{A}(\mathbf{x})) > \alpha\} = P_b\{\mathbf{x} \in \mathbb{R}^d : \widehat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) < q\}$$

and

$$\begin{aligned}
&P_b\{\mathbf{x} \in \mathbb{R}^d : P_\theta(\widehat{A}(\mathbf{x})) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d : P_\theta(A(\mathbf{x})) > \alpha\} \\
&= P_b\{\mathbf{x} \in \mathbb{R}^d : \widehat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) < q\} \\
&\quad - P_b\{\mathbf{x} \in \mathbb{R}^d : \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) < q\} \\
&= P_b\{\mathbf{x} \in \mathbb{R}^d : \widehat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_b) < \widehat{\mathbf{u}}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q\} \quad (3.18) \\
&\quad - P_b\{\mathbf{x} \in \mathbb{R}^d : \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_b) < \mathbf{u}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q\}.
\end{aligned}$$

Next we consider the second summand of (3.16). For $y \in \mathbb{R}$ let

$$\begin{aligned}
G(y) &:= P_\theta\{\mathbf{z} \in \mathbb{R}^d : \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{z} - \boldsymbol{\mu}_\theta) \geq y\} \\
&= P_0\{\mathbf{z} \in \mathbb{R}^d : \mathbf{u}^\top \mathbf{z} \geq y\} \\
&= 1 - F_{Z_1}(y)
\end{aligned}$$

and therefore $G'(y) = -f_{Z_1}(y) < 0$.

Note that

$$\begin{aligned} P_\theta(A(\mathbf{x})) &= P_\theta\{\mathbf{z} \in \mathbb{R}^d: \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{z} - \boldsymbol{\mu}_\theta) \geq \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)\} \\ &= G(\mathbf{u}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)) \end{aligned}$$

and the spherically symmetric distribution of $\Sigma^{-1/2}(\mathbf{X}_\theta - \boldsymbol{\mu}_\theta)$ implies

$$\begin{aligned} P_\theta(\hat{A}(\mathbf{x})) &= P_\theta\{\mathbf{z} \in \mathbb{R}^d: \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{z} - \boldsymbol{\mu}_\theta) \geq \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)\} \\ &= P_\theta\{\mathbf{z} \in \mathbb{R}^d: \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{z} - \boldsymbol{\mu}_\theta) \geq \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)\} \\ &= G(\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)). \end{aligned}$$

The first order Taylor expansion of G at q is given by

$$\begin{aligned} P_\theta(\hat{A}(\mathbf{x})) &= G(q) + G'(q) \cdot (\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q) \\ &\quad + o_p(\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q) \\ &= \alpha + G'(q) \cdot (\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q) \\ &\quad + o_p(\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q). \end{aligned}$$

For $\boldsymbol{\beta} \in \mathbb{R}^d$ and $\gamma \in \mathbb{R}$ consider the half-space $H(\boldsymbol{\beta}, \gamma) := \{\mathbf{z} \in \mathbb{R}^d: \boldsymbol{\beta}^\top \mathbf{z} + \gamma \leq 0\}$ and the empirical process $\mathbb{B}_{P_\theta, n}(\boldsymbol{\beta}, \gamma) := \sqrt{n}(\hat{P}_\theta - P_\theta)(H(\boldsymbol{\beta}, \gamma))$. Then

$$\begin{aligned} \hat{P}_\theta(\hat{A}(\mathbf{x})) &= P_\theta(\hat{A}(\mathbf{x})) + (\hat{P}_\theta - P_\theta)(\hat{A}(\mathbf{x})) \\ &= P_\theta(\hat{A}(\mathbf{x})) + n^{-1/2} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \hat{\mathbf{u}}, \hat{\mathbf{u}}^\top \Sigma^{-1/2} \mathbf{x}) \\ &= \alpha + G'(q) \cdot (\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q) \\ &\quad + n^{-1/2} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \hat{\mathbf{u}}, \hat{\mathbf{u}}^\top \Sigma^{-1/2} \mathbf{x}) \\ &\quad + o_p(\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q). \end{aligned}$$

Next we show that for the computation of

$$P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{P}_\theta(\hat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(\hat{A}(\mathbf{x})) > \alpha\}$$

it suffices to consider all \mathbf{x} such that

$$|\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q| < n^{-1/2} c^{-1} \|\mathbb{B}_{P_\theta, n}\|_\infty + O_p(n^{-1}) = O_p(n^{-1/2}) \quad (3.19)$$

for some constant $c > 0$. To this end note that

$$\mathbb{1}\{\hat{P}_\theta(\hat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} \neq \mathbb{1}\{P_\theta(\hat{A}(\mathbf{x})) > \alpha\}$$

3. Central Limit Theorems

implies

$$\begin{aligned}\|\mathbb{B}_{P_\theta, n}\|_\infty &> \sqrt{n}|P_\theta(\hat{A}(\mathbf{x})) - \alpha| + \mathcal{O}_p(n^{-1/2}) \\ &= \sqrt{n}|G(\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)) - \alpha| + \mathcal{O}_p(n^{-1/2}).\end{aligned}$$

Now suppose that $|\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q| \leq \delta$. Because $|G(q) - \alpha| = 0$,

$$|G(\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)) - \alpha| \geq |\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q| \min_{t \in [q \pm \delta]} |G'(t)|,$$

which entails (3.19) with $c := \min_{t \in [q \pm \delta]} |G'(t)|$.

If $|\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q| > \delta$, the monotonicity of G implies that

$$|G(\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta)) - \alpha| \geq \delta c$$

and hence $\|\mathbb{B}_{P_\theta, n}\|_\infty > \sqrt{n}\delta c + \mathcal{O}_p(n^{-1/2})$. But $\|\mathbb{B}_{P_\theta, n}\|_\infty = \mathcal{O}_p(1)$ by Theorem 3.4 and thus $\mathbb{P}(\|\mathbb{B}_{P_\theta, n}\|_\infty > \sqrt{n}\delta c + \mathcal{O}_p(n^{-1/2})) \rightarrow 0$.

Suppose that (3.19) holds. Then

$$P_\theta(\hat{A}(\mathbf{x})) = \alpha + G'(q) \cdot (\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q) + \mathcal{O}_p(n^{-1/2}).$$

Moreover, $\hat{\mathbf{u}} = \mathbf{u} + \mathcal{O}_p(n^{-1/2})$ and

$$\begin{aligned}\hat{\mathbf{u}}^\top \Sigma^{-1/2} \mathbf{x} &= \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) + \hat{\mathbf{u}}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta \\ &= q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta + \mathcal{O}_p(n^{-1/2}).\end{aligned}$$

The spherical symmetry of $\mathcal{L}(\mathbf{Z})$ and $\mathbb{E}(\|\mathbf{Z}\|^2) < \infty$ imply (3.2). Therefore the assumptions of Lemma 3.1 are satisfied. Applying the lemma, we get

$$\begin{aligned}P_\theta(H(-\Sigma^{-1/2}\hat{\mathbf{u}}, \hat{\mathbf{u}}^\top \Sigma^{-1/2} \mathbf{x}) \triangle H(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta)) \\ = \mathcal{O}_p(n^{-1/2})\end{aligned}$$

and by Theorem 3.5,

$$\mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2}\hat{\mathbf{u}}, \hat{\mathbf{u}}^\top \Sigma^{-1/2} \mathbf{x}) = \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + \mathcal{O}_p(1).$$

Consequently,

$$\begin{aligned}\hat{P}_\theta(\hat{A}(\mathbf{x})) &= \alpha + G'(q) \cdot (\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q) \\ &\quad + n^{-1/2} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + \mathcal{O}_p(n^{-1/2}).\end{aligned}$$

Moreover, $G'(q) < 0$ yields

$$\begin{aligned}
 & P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{P}_\theta(\hat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(\hat{A}(\mathbf{x})) > \alpha\} \\
 &= P_b\{\mathbf{x} \in \mathbb{R}^d: G'(q) \cdot (\hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_\theta) - q) \\
 &\quad + n^{-1/2} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + o_p(n^{-1/2}) > 0\} \\
 &\quad - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(\hat{A}(\mathbf{x})) > \alpha\} \\
 &= P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{\mathbf{u}}^\top \Sigma^{-1/2} \mathbf{x} < \hat{\mathbf{u}}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta + q \\
 &\quad - \frac{1}{\sqrt{n}G'(q)} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + o_p(n^{-1/2})\} \\
 &\quad - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(\hat{A}(\mathbf{x})) > \alpha\} \\
 &= P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_b) < \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q \quad (3.20) \\
 &\quad - \frac{1}{\sqrt{n}G'(q)} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + o_p(n^{-1/2})\} \\
 &\quad - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(\hat{A}(\mathbf{x})) > \alpha\}.
 \end{aligned}$$

Combining (3.18) and (3.20) we get

$$\begin{aligned}
 & P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{P}_\theta(\hat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(A(\mathbf{x})) > \alpha\} \\
 &= P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_b) < \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q \\
 &\quad - \frac{1}{\sqrt{n}G'(q)} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + o_p(n^{-1/2})\} \\
 &\quad - P_b\{\mathbf{x} \in \mathbb{R}^d: \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{x} - \boldsymbol{\mu}_b) < \mathbf{u}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q\} \\
 &= P_0\{\mathbf{x} \in \mathbb{R}^d: \hat{\mathbf{u}}^\top \mathbf{x} < \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q \\
 &\quad - \frac{1}{\sqrt{n}G'(q)} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + o_p(n^{-1/2})\} \\
 &\quad - P_0\{\mathbf{x} \in \mathbb{R}^d: \mathbf{u}^\top \mathbf{x} < \mathbf{u}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q\}
 \end{aligned}$$

with $P_0 := \mathcal{L}(\Sigma^{-1/2}(\mathbf{X}_b - \boldsymbol{\mu}_b))$. The spherical symmetry of P_0 yields

$$\begin{aligned}
 & P_b\{\mathbf{x} \in \mathbb{R}^d: \hat{P}_\theta(\hat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(A(\mathbf{x})) > \alpha\} \\
 &= P_0\{\mathbf{x} \in \mathbb{R}^d: x_1 < \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q \\
 &\quad - \frac{1}{\sqrt{n}G'(q)} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + o_p(n^{-1/2})\} \\
 &\quad - P_0\{\mathbf{x} \in \mathbb{R}^d: x_1 < \mathbf{u}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q\} \\
 &= P_0\{\mathbf{x} \in \mathbb{R}^d: x_1 - \hat{s} \leq 0\} - P_0\{\mathbf{x} \in \mathbb{R}^d: x_1 - s \leq 0\}
 \end{aligned}$$

with $\hat{s} := \hat{\mathbf{u}}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q - \frac{1}{\sqrt{n}G'(q)} \mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2} \mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2} \boldsymbol{\mu}_\theta) + o_p(n^{-1/2})$ and $s := \mathbf{u}^\top \Sigma^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) + q = q - \|\Sigma^{-1/2}(\boldsymbol{\mu}_b - \boldsymbol{\mu}_\theta)\|$.

3. Central Limit Theorems

Since $\widehat{s} - s = O_p(n^{-1/2})$ we can apply Lemma 3.1 with $\widehat{\mathbf{b}} = \mathbf{b} = \mathbf{e}_1$, $\widehat{a} = -\widehat{s}$ and $a = -s$. We obtain claim (3.14) and

$$\begin{aligned} P_b\{\mathbf{x} \in \mathbb{R}^d: \widehat{P}_\theta(\widehat{A}(\mathbf{x})) + O_p(n^{-1}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: P_\theta(A(\mathbf{x})) > \alpha\} \\ = (\widehat{s} - s)f_{Z_1}(s) + O_p(n^{-1}). \end{aligned}$$

If $b = \theta$, $(\widehat{\mathbf{u}} - \mathbf{u})^\top \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\mu}_\theta - \boldsymbol{\mu}_b) = 0$. Suppose that $b \neq \theta$. The first order Taylor expansion of $\|\widehat{\boldsymbol{\beta}}\|$ is then given by

$$\|\widehat{\boldsymbol{\beta}}\| = \|\boldsymbol{\beta}\| + \frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \boldsymbol{\beta}}{\|\boldsymbol{\beta}\|} + o(\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|)$$

and therefore

$$\begin{aligned} \widehat{\mathbf{u}} - \mathbf{u} &= \frac{\|\boldsymbol{\beta}\|\widehat{\boldsymbol{\beta}} - \|\widehat{\boldsymbol{\beta}}\|\boldsymbol{\beta}}{\|\boldsymbol{\beta}\|\|\widehat{\boldsymbol{\beta}}\|} \\ &= \frac{\|\boldsymbol{\beta}\|(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) - \frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \boldsymbol{\beta}}{\|\boldsymbol{\beta}\|}\boldsymbol{\beta}}{\|\boldsymbol{\beta}\|^2 + (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \boldsymbol{\beta}} + o(\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|) \\ &= \frac{\|\boldsymbol{\beta}\|(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) - \frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \boldsymbol{\beta}}{\|\boldsymbol{\beta}\|}\boldsymbol{\beta}}{\|\boldsymbol{\beta}\|^2} + o(\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|) \\ &= \frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{\|\boldsymbol{\beta}\|} - \frac{(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \boldsymbol{\beta}}{\|\boldsymbol{\beta}\|^3}\boldsymbol{\beta} + o(\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|). \end{aligned}$$

Plugging in $\beta = \Sigma^{-1/2}(\mu_\lambda - \mu_\theta)$ and $\hat{\beta} := \beta + \Sigma^{1/2} \frac{\Delta_{\Sigma^{-1}}}{\sqrt{n}}(\mu_\lambda - \mu_\theta) + \Sigma^{-1/2} \frac{\Delta_{\mu_\lambda} - \Delta_{\mu_\theta}}{\sqrt{n}} + O_p(n^{-1})$ results in

$$\begin{aligned} \hat{u} - u &= \frac{\Sigma^{1/2} \Delta_{\Sigma^{-1}}(\mu_b - \mu_\theta) + \Sigma^{-1/2}(\Delta_{\mu_b} - \Delta_{\mu_\theta})}{\sqrt{n} \|\Sigma^{-1/2}(\mu_b - \mu_\theta)\|} \\ &\quad - \frac{(\mu_b - \mu_\theta)^\top \Delta_{\Sigma^{-1}}(\mu_b - \mu_\theta) + (\Delta_{\mu_b} - \Delta_{\mu_\theta})^\top \Sigma^{-1}(\mu_b - \mu_\theta)}{\sqrt{n} \|\Sigma^{-1/2}(\mu_b - \mu_\theta)\|^3} \\ &\quad \cdot \Sigma^{-1/2}(\mu_b - \mu_\theta) + o_p(n^{-1/2}). \end{aligned}$$

Hence

$$\begin{aligned} &(\hat{u} - u)^\top \Sigma^{-1/2}(\mu_\theta - \mu_b) \\ &= \frac{(\mu_b - \mu_\theta)^\top \Delta_{\Sigma^{-1}}(\mu_\theta - \mu_b) + (\Delta_{\mu_b} - \Delta_{\mu_\theta})^\top \Sigma^{-1}(\mu_\theta - \mu_b)}{\sqrt{n} \|\Sigma^{-1/2}(\mu_b - \mu_\theta)\|} \\ &\quad + \frac{(\mu_b - \mu_\theta)^\top \Delta_{\Sigma^{-1}}(\mu_b - \mu_\theta) + (\Delta_{\mu_b} - \Delta_{\mu_\theta})^\top \Sigma^{-1}(\mu_b - \mu_\theta)}{\sqrt{n} \|\Sigma^{-1/2}(\mu_b - \mu_\theta)\|} \\ &\quad + o_p(n^{-1/2}) \\ &= o_p(n^{-1/2}) \end{aligned}$$

and

$$\begin{aligned} &P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) > \alpha\} \\ &= (\hat{s} - s)f_{Z_1}(s) + O_p(n^{-1}) \\ &= -\frac{1}{\sqrt{n}G'(q)}\mathbb{B}_{P_\theta, n}(-\Sigma^{-1/2}\mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2}\mu_\theta)f_{Z_1}(s) + o_p(n^{-1/2}). \end{aligned}$$

Finally, with

$$\begin{aligned} H(-\Sigma^{-1/2}\mathbf{u}, q + \mathbf{u}^\top \Sigma^{-1/2}\mu_\theta) &= \{\mathbf{x} \in \mathbb{R}^d: \mathbf{u}^\top \Sigma^{-1/2}(\mathbf{x} - \mu_\theta) \geq q\} \\ &= \{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) \leq \alpha\} \end{aligned}$$

we get

$$\begin{aligned} &P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\} - P_b\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) > \alpha\} \\ &= \frac{f_{Z_1}(q - \|\Sigma^{-1/2}(\mu_b - \mu_\theta)\|)}{f_{Z_1}(q)}(\hat{P}_\theta - P_\theta)(\{\mathbf{x} \in \mathbb{R}^d: \pi_\theta^*(\mathbf{x}) \leq \alpha\}) \\ &\quad + o_p(n^{-1/2}). \end{aligned}$$

□

3. Central Limit Theorems

PROOF OF THEOREM 3.11. Define $\widehat{E} := \{\mathbf{x} \in \mathbb{R}^d : \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\}$ and note that $\widehat{\mathcal{I}}_\alpha(b, \theta) = \widehat{P}_b(\widehat{E}) + o_p(n^{-1/2})$. Then by Theorem 3.5 and equation (3.14) from Lemma 3.13,

$$\begin{aligned} \widehat{\mathcal{I}}_\alpha(b, \theta) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = b) \\ &= \widehat{P}_b(\widehat{E}) - P_b(E) + o_p(n^{-1/2}) \\ &= (P_b(\widehat{E}) - P_b(E)) + (\widehat{P}_b - P_b)(E) \\ &\quad + \left((\widehat{P}_b - P_b)(\widehat{E}) - (\widehat{P}_b - P_b)(E) \right) + o_p(n^{-1/2}) \\ &= (P_b(\widehat{E}) - P_b(E)) + (\widehat{P}_b - P_b)(E) + o_p(n^{-1/2}). \end{aligned}$$

Let

$$\begin{aligned} \boldsymbol{\eta} &:= \sqrt{n} \begin{pmatrix} \mathcal{I}_\alpha(b, \theta \mid \mathcal{D}) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = b) \\ \widehat{\mathcal{I}}_\alpha(b, \theta) - \mathbb{P}(\pi_\theta^*(\mathbf{X}) > \alpha \mid Y = b) \end{pmatrix} \\ &= \sqrt{n} \begin{pmatrix} P_b(\widehat{E}) - P_b(E) \\ (P_b(\widehat{E}) - P_b(E)) + (\widehat{P}_b - P_b)(E) \end{pmatrix} + o_p(1). \end{aligned}$$

We employ Lemma 3.13 to decompose $P_b(\widehat{E}) - P_b(E)$ in independent summands

$$\begin{aligned} P_b(\widehat{E}) - P_b(E) &= -c(\widehat{P}_\theta - P_\theta)\{\mathbf{x} \in \mathbb{R}^d : \pi_\theta^*(\mathbf{x}) \leq \alpha\} + o_p(n^{-1/2}) \\ &= c(\widehat{P}_\theta - P_\theta)(E) + o_p(n^{-1/2}), \end{aligned}$$

where $c = -f_{Z_1}(q - \|\boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\mu}_b - \boldsymbol{\mu}_\theta)\|)/f_{Z_1}(q)$. Note that $c = -1$ if $b = \theta$. In this case

$$\begin{aligned} \widehat{P}_\theta\{\mathbf{x} \in \mathbb{R}^d : \pi_\theta(\mathbf{x}, \mathcal{D}) > \alpha\} - P_\theta\{\mathbf{x} \in \mathbb{R}^d : \pi_\theta^*(\mathbf{x}) > \alpha\} \\ &= -(\widehat{P}_\theta - P_\theta)(E) + (\widehat{P}_\theta - P_\theta)(E) + o_p(n^{-1/2}) \\ &= o_p(n^{-1/2}), \end{aligned}$$

which is claim (3.13).

Now we define for $\nu \in \{b, \theta\}$

$$Y_{n,i}^\nu := \frac{1}{\sqrt{n\widehat{w}_\nu}} \mathbf{1}\{i \in \mathcal{G}_\nu\} (\mathbf{1}\{\mathbf{X}_i \in E\} - P_\nu(E))$$

such that

$$\begin{aligned} \sum_{i=1}^n Y_{n,i}^\nu &= \frac{\sqrt{n}}{N_\nu} \sum_{i \in \mathcal{G}_\nu} (\mathbf{1}\{\mathbf{X}_i \in E\} - P_\nu(E)) \\ &= \sqrt{n}(\widehat{P}_\nu - P_\nu)(E). \end{aligned}$$

Thus we may write

$$\boldsymbol{\eta} = \sum_{i=1}^n \mathbf{Y}_{n,i}^{\boldsymbol{\eta}} + o_p(1)$$

with

$$\mathbf{Y}_{n,i}^{\boldsymbol{\eta}} := \begin{pmatrix} cY_{n,i}^{\theta} \\ cY_{n,i}^{\theta} + Y_{n,i}^b \end{pmatrix}.$$

Before we can apply the central limit theorem to $\sum_{i=1}^n \mathbf{Y}_{n,i}^{\boldsymbol{\eta}}$ we need to compute the covariance matrix $\text{Var}(\mathbf{Y}_{n,i}^{\boldsymbol{\eta}})$, the sum $\boldsymbol{\Lambda}^n := \sum_{i=1}^n \text{Var}(\mathbf{Y}_{n,i}^{\boldsymbol{\eta}})$ and its limit $\boldsymbol{\Lambda} = \lim_{n \rightarrow \infty} \boldsymbol{\Lambda}^n$. To this end note that

$$\text{Var}(Y_{n,i}^{\nu}) = \frac{\mathbb{1}\{i \in \mathcal{G}_{\nu}\}}{n\widehat{w}_{\nu}^2} (P_{\nu}(E) - P_{\nu}(E)^2)$$

and

$$\begin{aligned} \Lambda_{1,1}^n &= c^2 \sum_{i=1}^n \text{Var}(Y_{n,i}^{\theta}) \\ &\rightarrow \frac{c^2}{w_{\theta}} (P_{\theta}(E) - P_{\theta}(E)^2) \\ &= w_{\theta}^{-1} c^2 \alpha (1 - \alpha). \end{aligned}$$

Suppose that $b \neq \theta$. Then $\text{Cov}(Y_{n,i}^{\theta}, Y_{n,i}^b) = 0$ and thus $\Lambda_{1,1} = \Lambda_{1,2} = \Lambda_{2,1}$. Moreover,

$$\begin{aligned} \text{Var}(cY_{n,i}^{\theta} + Y_{n,i}^{\lambda}) &= c^2 \text{Var}(Y_{n,i}^{\theta}) + \text{Var}(Y_{n,i}^{\lambda}) \\ &= \frac{c^2 \mathbb{1}\{i \in \mathcal{G}_{\theta}\}}{n\widehat{w}_{\theta}^2} (P_{\theta}(E) - P_{\theta}(E)^2) \\ &\quad + \frac{\mathbb{1}\{i \in \mathcal{G}_{\lambda}\}}{n\widehat{w}_{\lambda}^2} (P_{\lambda}(E) - P_{\lambda}(E)^2) \end{aligned}$$

and

$$\begin{aligned} \Lambda_{2,2}^n &= c^2 \widehat{w}_{\theta}^{-1} (P_{\theta}(E) - P_{\theta}(E)^2) + \widehat{w}_{\lambda}^{-1} (P_{\lambda}(E) - P_{\lambda}(E)^2) \\ &\rightarrow w_{\theta}^{-1} c^2 \alpha (1 - \alpha) + w_{\lambda}^{-1} P_{\lambda}(E) (1 - P_{\lambda}(E)) \end{aligned}$$

with

$$\begin{aligned} P_{\lambda}(E) &= P_0\{\mathbf{x} \in \mathbb{R}^d: \mathbf{u}^{\top} \mathbf{x} < q + \mathbf{u}^{\top} \boldsymbol{\Sigma}^{-1/2} (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_{\lambda})\} \\ &= P_0\{\mathbf{x} \in \mathbb{R}^d: \mathbf{u}^{\top} \mathbf{x} < q - \|\boldsymbol{\beta}\|\} \\ &= F_{Z_1}(q - \|\boldsymbol{\beta}\|), \end{aligned}$$

3. Central Limit Theorems

which follows from equation (3.17).

As to Lindeberg's condition, note that $|Y_{n,i}^\nu| \leq 2n^{-1/2}\widehat{w}_\nu^{-1}$ and hence

$$\|\mathbf{Y}_{n,i}^\eta\| \leq 2|c||Y_{n,i}^\theta| + |Y_{n,i}^b| \leq \frac{4|c|}{\sqrt{n}\widehat{w}_\theta} + \frac{2}{\sqrt{n}\widehat{w}_b} \leq \frac{M_n}{\sqrt{n}},$$

where M_n is deterministic and bounded. Therefore Lindeberg's condition

$$\sum_{i=1}^n \mathbb{E}(\|\mathbf{Y}_{n,i}^\eta\|^2 \min(1, \|\mathbf{Y}_{n,i}^\eta\|)) \leq n^{-1/2} M_n^3 \rightarrow 0$$

is satisfied. Finally note that $\mathbb{E}(\mathbf{Y}_{n,i}^\eta) = \mathbf{0}$ for all $i \leq n$. Now the assertions follow from the multivariate central limit theorem (Theorem A.1). \square

4. Randomized and De-Randomized P-Values

In this chapter we discuss the concept of randomized p-values which is not directly related to the main topic of this thesis. The corresponding concept of randomized tests is familiar in mathematical statistics and is used to obtain tests with exact prescribed significance levels even in settings with test statistics having discrete distributions. Similarly, randomized p-values are particularly useful for test statistics with discrete distributions.

In applications non-randomized tests and p-values are needed. Therefore we review and modify a method of Meinshausen et al. (2009) to de-randomize p-values. To randomize and de-randomize one single p-value brings no benefit, since the resulting p-value is greater than the initial one. However, if we want to combine p-values obtained from different independent test statistics, randomization may be useful. Randomize the p-values, combine them and de-randomize the combination again can lead to a result which is considerably smaller than the combination of the non-randomized p-values.

4.1. De-Randomization

Example 4.1. We consider a random variable $X \sim \text{Poiss}(\lambda)$ with unknown parameter λ and want to test $H_0: \lambda = \lambda_0$ versus $H_A: \lambda > \lambda_0$. The usual p-value is given by $1 - G_{\lambda_0}(T - 1)$ with the test statistic $T = X$ and G_{λ_0} denoting the c.d.f. of $\text{Poiss}(\lambda_0)$. This p-value is conservative due to the discrete distribution of T .

In mathematical statistics we look at the randomized level- α -test rejecting H_0 with probability

$$\varphi(x) = \begin{cases} 1 & \text{if } 1 - G_{\lambda_0}(T - 1) \leq \alpha \\ \gamma & \text{if } 1 - G_{\lambda_0}(T) \leq \alpha < 1 - G_{\lambda_0}(T - 1) \\ 0 & \text{if } \alpha < 1 - G_{\lambda_0}(T), \end{cases}$$

where

$$\gamma := \frac{\mathbb{P}(1 - G_{\lambda_0}(T - 1)) - \alpha}{\mathbb{P}(1 - G_{\lambda_0}(T - 1)) - \mathbb{P}(1 - G_{\lambda_0}(T))}.$$

4. Randomized and De-Randomized P-Values

The corresponding randomized p-value is given by

$$\pi \sim \text{Unif}[1 - G_{\lambda_0}(T), 1 - G_{\lambda_0}(T - 1)].$$

In applications non-randomized tests and p-values are needed. Therefore we will de-randomize this p-value.

The abstract setting. Consider a probability space $(\Omega, \mathcal{A}, \mathbb{P})$, and let $T : (\Omega, \mathcal{A}) \rightarrow (\mathcal{T}, \mathcal{B})$ and $\pi : (\Omega, \mathcal{A}) \rightarrow [0, 1]$ be measurable mappings satisfying the following assumptions under a certain null hypothesis:

(A.1) $\mathcal{L}(\pi) = \text{Unif}[0, 1]$.

(A.2) For each $t \in \mathcal{T}$ there exists a given distribution function F_t on $[0, 1]$ such that

$$\mathbb{P}(\pi \leq u | T) = F_T(u) \quad \text{almost surely, for each } u \in [0, 1].$$

Note that π is a randomized p-value, and the pair (T, π) may be represented as

$$(T, \pi) = (T, Q_T(U)),$$

where T and U are independent, $U \sim \text{Unif}[0, 1]$, and Q_t is the quantile function of F_t , i.e.

$$Q_t(v) := \min\{u \in [0, 1] : F_t(u) \geq v\}.$$

In our specific applications, T is a single or vector-valued test statistic, and it is often desirable to come up with a p-value $\tilde{\pi}$ depending on T only. A naive solution would be

$$\tilde{\pi} := Q_T(1),$$

because $Q_T(1) \geq \pi$ almost surely. But this may be much too conservative. Here is a first general proposal how to construct $\tilde{\pi}$:

De-randomization in the spirit of Meinshausen et al. (2009). Let $\Gamma \subset (0, 1]$ be a nonvoid set, and let $h : \Gamma \rightarrow (0, \infty)$. Defining

$$J(u) := \sup_{\gamma \in \Gamma} \frac{\mathbb{1}\{u \leq \gamma\}}{h(\gamma)} \quad \text{for } u \geq 0,$$

we assume that

$$J := \int_0^1 J(u) du < \infty.$$

We define $Q_t(v) := \infty$ for $v > 1$ and note that $J(u) = 0$ for $u > 1$. Then for any $\alpha \in (0, 1)$,

$$\begin{aligned}
\alpha &= \mathbb{E}\left(\frac{J(\pi/\alpha)}{J}\right) \\
&= \mathbb{E}\left(\mathbb{E}\left(\frac{J(\pi/\alpha)}{J} \mid T\right)\right) \\
&\geq \mathbb{E}\left(\sup_{\gamma \in \Gamma} \mathbb{E}\left(\frac{\mathbf{1}\{\pi \leq \alpha\gamma\}}{Jh(\gamma)} \mid T\right)\right) \\
&= \mathbb{E}\left(\sup_{\gamma \in \Gamma} \frac{F_T(\alpha\gamma)}{Jh(\gamma)}\right) \\
&\geq \mathbb{E}\left(\sup_{\gamma \in \Gamma} \mathbf{1}\{F_T(\alpha\gamma) \geq Jh(\gamma)\}\right) \\
&= \mathbb{E}\left(\sup_{\gamma \in \Gamma} \mathbf{1}\{Q_T(Jh(\gamma)) \leq \alpha\gamma\}\right) \\
&= \mathbb{P}(Q_T(Jh(\gamma)) \leq \alpha\gamma \text{ for some } \gamma \in \Gamma).
\end{aligned}$$

Thus we may reject the null hypothesis at level α if

$$\frac{Q_T(Jh(\gamma))}{\gamma} \leq \alpha \text{ for some } \gamma \in \Gamma.$$

A corresponding non-randomized p-value is given by

$$\tilde{\pi} := \min_{\gamma \in \Gamma} \frac{Q_T(Jh(\gamma))}{\gamma},$$

provided the latter minimum exists almost surely.

Example 4.2. Let $\Gamma = \{\gamma_o\}$ for some fixed $\gamma_o \in (0, 1)$ and $h(\gamma_o) = \gamma_o$. Then $J(u) = \mathbf{1}\{u \leq \gamma_o\}/\gamma_o$, $J = 1$, and the nonrandomized p-value is given by

$$\tilde{\pi} := \frac{Q_T(\gamma_o)}{\gamma_o}.$$

Example 4.3. Let $\Gamma = [\gamma_o, 1]$ for some fixed $\gamma_o \in (0, 1)$ and $h(\gamma) = \gamma$. Then

$$J(u) = \sup_{\gamma_o \leq \gamma \leq 1} \frac{\mathbf{1}\{u \leq \gamma\}}{\gamma} = \begin{cases} 1/\gamma_o & \text{if } 0 \leq u \leq \gamma_o, \\ 1/u & \text{if } \gamma_o \leq u \leq 1. \end{cases}$$

Thus

$$J = \int_0^{\gamma_o} \frac{1}{\gamma_o} du + \int_{\gamma_o}^1 \frac{1}{u} du = \log(e/\gamma_o).$$

4. Randomized and De-Randomized P-Values

Consequently, a nonrandomized p-value is given by

$$\begin{aligned}\tilde{\pi} &:= \min_{\gamma_o \leq \gamma \leq 1} \frac{Q_T(\log(e/\gamma_o)\gamma)}{\gamma} \\ &= \min_{\gamma_o \leq \gamma \leq 1/\log(e/\gamma_o)} \frac{Q_T(\log(e/\gamma_o)\gamma)}{\gamma} \\ &= \min_{\gamma_o \log(e/\gamma_o) \leq u \leq 1} \frac{\log(e/\gamma_o)Q_T(u)}{u}.\end{aligned}$$

Example 4.4. Let $\Gamma = (0, 1]$ and $h(\gamma) = \gamma^\delta$ for some $\delta \in (0, 1)$. Then

$$J(u) = \sup_{0 < \gamma \leq 1} \frac{\mathbb{1}\{u \leq \gamma\}}{\gamma^\delta} = \mathbb{1}\{u \leq 1\}u^{-\delta} \quad \text{and} \quad J = (1 - \delta)^{-1}.$$

Consequently, a nonrandomized test rejects the null hypothesis if

$$\frac{Q_T(\gamma^\delta/(1 - \delta))}{\gamma} \leq \alpha \quad \text{for some } \gamma \in (0, 1],$$

which is equivalent to

$$\frac{Q_T(u)}{((1 - \delta)u)^{1/\delta}} \leq \alpha \quad \text{for some } u \in (0, 1].$$

The corresponding p-value,

$$\tilde{\pi} = \min_{0 < u \leq 1} \frac{Q_T(u)}{((1 - \delta)u)^{1/\delta}},$$

is well-defined if, for instance,

$$\limsup_{u \downarrow 0} \frac{F_t(u)}{u} < \infty \quad \text{for all } t \in \mathcal{T}.$$

For then, $Q_t(u) \geq c(t)u$ for all $u \in (0, 1]$ and some $c(t) > 0$, so that $Q_t(u)/u^{1/\delta} \rightarrow \infty$ as $u \downarrow 0$. Moreover, Q_t is left-continuous and non-decreasing, and this entails that $Q_t(u)/u^{1/\delta}$ attains a minimum on $(0, 1]$.

4.2. Combining Independent P-Values

Suppose that for a given null hypothesis, stochastically independent and possibly randomized p-values $\pi_1, \pi_2, \dots, \pi_m$ are available. There are infinitely

many possibilities to combine these p-values into one p-value. One specific way is to use

$$\pi := \Phi\left(m^{-1/2} \sum_{i=1}^m \Phi^{-1}(\pi_i)\right) \quad (4.1)$$

with the standard Gaussian distribution function Φ . More generally, we may define

$$\pi := \Phi\left(\sum_{i=1}^m w_i \Phi^{-1}(\pi_i)\right) \quad (4.2)$$

with certain weights $w_1, w_2, \dots, w_m > 0$ such that $\sum_{i=1}^m w_i^2 = 1$.

This method is motivated by the following model: Suppose we observe independent random variables X_1, X_2, \dots, X_m with $X_i \sim \mathcal{N}(\theta, \sigma_i^2)$, where the standard deviations $\sigma_i > 0$ are known while the mean θ is unknown. For the null hypothesis “ $\theta \geq 0$ ”, possible p-values are given by

$$\pi_i := \Phi\left(\frac{X_i}{\sigma_i}\right).$$

Then proposal (4.1) leads to

$$\Phi\left(m^{-1/2} \sum_{i=1}^m \frac{X_i}{\sigma_i}\right) \sim \text{Unif}[0, 1] \quad \text{if } \theta = 0.$$

In this simple model, the Neyman-Pearson Lemma (Theorem A.2) shows that an optimal p-value is given by

$$\Phi\left(\left(\sum_{i=1}^m \sigma_i^{-2}\right)^{-1/2} \sum_{i=1}^m \frac{X_i}{\sigma_i^2}\right) = \Phi\left(\left(\sum_{i=1}^m \sigma_i^{-2}\right)^{-1/2} \sum_{i=1}^m \frac{\Phi^{-1}(\pi_i)}{\sigma_i}\right),$$

and this corresponds to proposal (4.2) with

$$w_i := \left(\sum_{j=1}^m \sigma_j^{-2}\right)^{-1/2} \sigma_i^{-1}.$$

Here the optimal weights are proportional to the reciprocal standard deviations. But there is another representation which can be imitated in different settings: Note that with $Z \sim \mathcal{N}(0, 1)$,

$$\frac{d}{d\theta}\Big|_{\theta=0} \mathbb{E}_\theta \pi_i = \frac{d}{d\theta}\Big|_{\theta=0} \mathbb{E} \Phi\left(\frac{\theta + \sigma_i Z}{\sigma_i}\right) = \frac{1}{\sigma_i} \mathbb{E} \Phi'(Z) = \frac{\text{const.}}{\sigma_i}.$$

Thus in more general settings where our null hypothesis is of the form “ $\theta \geq 0$ ” for some real parameter θ , we propose to choose w_i proportional to

$$\frac{d}{d\theta}\Big|_{\theta=0} \mathbb{E}_\theta \pi_i.$$

4.3. Application to Multiple Contingency Tables

Suppose that we have multiple independent contingency tables, e.g. data from a multicenter clinical trial or a meta-analysis. Assume that the tables have small cell counts. In this case the test statistics can take only few values and randomization may be useful.

4.3.1. Two-by-Two Tables

Let $S^{(1)}, S^{(2)}, \dots, S^{(m)}$ be independent two-by-two tables:

X_i	$z_i - X_i$	z_i	
$s_i - X_i$	$n_i - s_i - z_i + X_i$	$n_i - z_i$	
s_i	$n_i - s_i$	n_i	

We consider the row and column sums as fixed and assume that the tables have a common but unknown odds ratio ρ .

Now we want to test the null hypothesis $\rho = 1$ versus the alternative $\rho < 1$. Under the null hypothesis X_i has a hypergeometric distribution with parameters n_i, z_i and s_i . A p-value for table $S^{(i)}$ is given by $\pi_i^{\text{cons}} := H_{n_i, z_i, s_i}(X_i)$ where $H_{n_i, z_i, s_i}(k) := \sum_{l=0}^k \binom{z_i}{l} \binom{n-z_i}{s_i-l} / \binom{n_i}{s_i}$ denotes the cdf of the hypergeometric distribution with parameters n_i, z_i and s_i . Due to its discrete distribution, this p-value may be rather conservative. Therefore we consider the randomized p-value

$$\pi_i \sim \text{Unif}[H_{n_i, z_i, s_i}(X_i - 1), H_{n_i, z_i, s_i}(X_i)].$$

Combining the p-values of all tables, we get the randomized p-value

$$\pi := \Phi\left(\sum_{i=1}^m w_i \Phi^{-1}(\pi_i)\right),$$

which we de-randomize as described in section 4.1. Note that the non-randomized p-values π_i^{cons} can be combined the same way, which results in the conservative p-value $\pi^{\text{cons}} := \Phi\left(\sum_{i=1}^m w_i \Phi^{-1}(\pi_i^{\text{cons}})\right)$.

In order to compute Q_T with $T = (X_i)_{i=1}^m$, we first approximate F_T numerically. To this end we discretize the distribution of the $Z_i := w_i \Phi^{-1}(\pi_i)$. For a fixed $\delta > 0$ we choose $C_{i,1}, C_{i,2} \in \mathbb{Z}\delta$ such that

$$\begin{aligned} \mathbb{P}(Z_i \leq C_{i,1} \mid T) &\ll 1, \\ \mathbb{P}(Z_i > C_{i,2} \mid T) &\ll 1. \end{aligned}$$

Then we define

$$\tilde{Z}_i := \begin{cases} C_{i,1} & \text{if } Z_i \leq C_{i,1}, \\ \lceil Z_i/\delta \rceil \delta & \text{if } C_{i,1} < Z_i < C_{i,2}, \\ C_{i,2} & \text{if } Z_i \geq C_{i,2}, \end{cases}$$

and compute $p_i = (p_i(j))_{j=1}^{M(i)}$, with $p_i(j) := \mathbb{P}(\tilde{Z}_i = C_{i,1} + (j-1)\delta \mid T)$. The approximate distribution of π conditional on T is then given by the convolution of p_1, p_2, \dots, p_m and its domain is given by

$$\Phi\left(\sum_{i=1}^m C_{i,1} + k\delta\right) \quad k = 0, 1, \dots, \sum_{i=1}^m (M(i) - 1),$$

where $M(i) := (C_{i,2} - C_{i,1})/\delta + 1$.

Simulation of X_i . To compute the power of the resulting test we need to simulate X_i under the alternative hypothesis. In general X_i has a non-central hypergeometric distribution, i.e.

$$\mathbb{P}(X_i = k \mid n_i, z_i, s_i) = f_\rho(k \mid n_i, z_i, s_i),$$

with

$$f_\rho(k \mid n, z, s) := C_\rho(n, z, s)^{-1} \rho^k \binom{z}{k} \binom{n-z}{s-k},$$

$$C_\rho(n, z, s) := \sum_{l=\max(z+s-n, 0)}^{\min(z, s)} \rho^l \binom{z}{l} \binom{n-z}{s-l},$$

see e.g. Agresti (2007). For the computation we use the representation

$$f_\rho(k \mid n, z, s) = \tilde{C}_\rho(n, z, s)^{-1} \frac{\rho^k}{k!(z-k)!(s-k)!(n-z-s+k)!},$$

$$\tilde{C}_\rho(n, z, s) := \sum_{l=\max(z+s-n, 0)}^{\min(z, s)} \frac{\rho^l}{l!(z-l)!(s-l)!(n-z-s+l)!}$$

and to avoid numerical problems, we compute it in three steps:

$$\begin{aligned} f &\leftarrow \log(\rho)k - \log(k!) - \log((z-k)!) - \log((s-k)!) \\ &\quad - \log((n-s-z+k)!) \\ f &\leftarrow \exp((f - \max(f))) \\ f &\leftarrow f/\text{sum}(f). \end{aligned}$$

4. Randomized and De-Randomized P-Values

	$\rho = 1$	$\rho = 0.8$	$\rho = 0.5$	$\rho = 1/3$
exact	0.0355	0.1443	0.7287	0.9809
conservative	0.0013	0.0147	0.2949	0.8189
$\beta = 0.1$	0.0064	0.0518	0.5204	0.9298
$\beta = 0.3$	0.0135	0.0907	0.6328	0.9615
$\beta = 0.5$	0.0135	0.0907	0.6328	0.9615
$\beta = 0.7$	0.0135	0.0907	0.6328	0.9615
$\beta = 0.9$	0.0135	0.0907	0.6328	0.9615
$\gamma_0 = 0.5$	0.0064	0.0518	0.5204	0.9298
$\gamma_0 = 0.1$	0.0031	0.0280	0.4064	0.8830

Table 4.1.: Power for Example 4.5 at significance level $\alpha = 0.05$.

Exact Monte-Carlo p-values. In the case of two-by-two tables we can compute exact Monte-Carlo p-values. We use them as a benchmark for the de-randomized p-values. In practice the de-randomized p-values are only useful for K -by- L tables with $\max(K, L) > 2$.

For the computation of the Monte-Carlo p-values note that

$$\log \mathbb{P}_\theta(X_i = k \mid n_i, z_i, s_i) = k\theta + \log(C_\rho(n, z, s)^{-1}) + \log\left(\binom{z}{k}\binom{n-z}{s-k}\right)$$

with $\theta := \log \rho$. The log-likelihood function of the whole model is given by

$$L(\theta) = \sum_{i=1}^m \left(\theta X_i + \log\left(\binom{z}{X_i}\binom{n-z}{s-X_i}\right) \right) + m \log(C_\rho(n, z, s)^{-1})$$

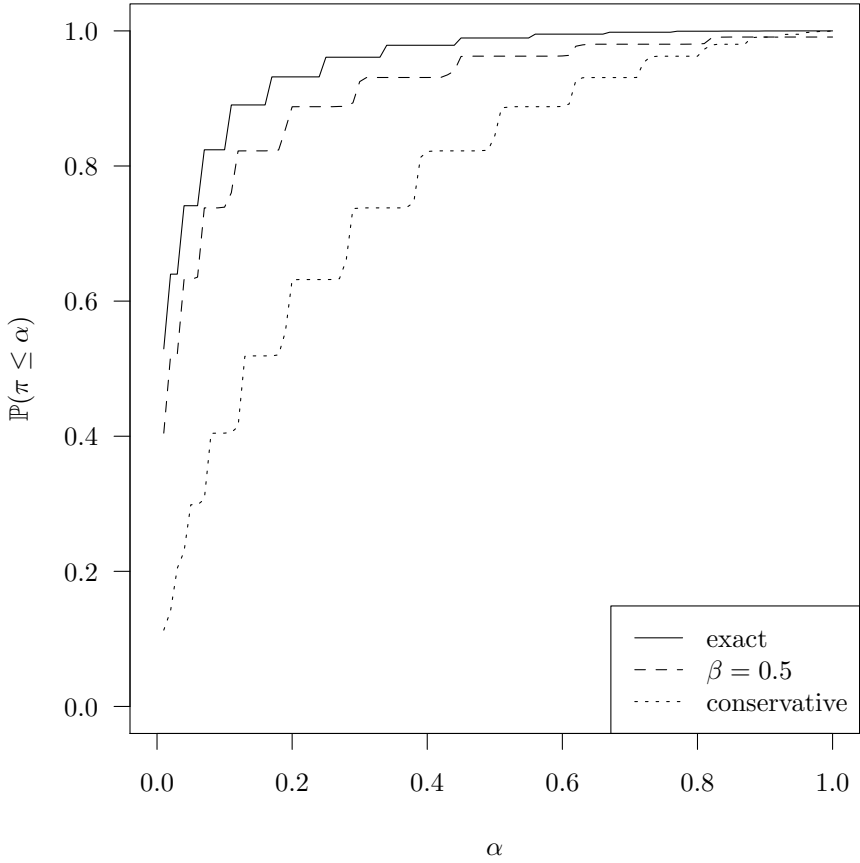
and a potential test statistic for " $\theta = \theta_0$ " vs. " $\theta < \theta_0$ " would be

$$\frac{\partial}{\partial \theta} \Big|_{\theta=\theta_0} L(\theta) = \sum_{i=1}^m X_i + m \frac{\partial}{\partial \theta} \Big|_{\theta=\theta_0} \log(C_\rho(n, z, s)^{-1}).$$

Since the last summand does not depend on the data, we choose the test statistic

$$\sum_{i=1}^m X_i.$$

Example 4.5. We simulated 10'000 times $m = 10$ tables with $n = 20$, $z = 10$, $s = 8$ and different values for the odds ratio ρ . To combine the p-values we used equal weights $w_i = m^{-1/2}$. Table 4.1 shows the power at significance level $\alpha = 0.05$ for the exact Monte-Carlo p-value, the conservative p-value π^{cons} , $Q_T(\beta)/\beta$ for $\beta = 0.1, 0.3, \dots, 0.9$ and the adaptive version

Figure 4.1.: Power for $\rho = 0.5$ in Example 4.5.

4. Randomized and De-Randomized P-Values

(Example 4.3) with $\gamma_0 = 0.5$ and 0.1 . Figure 4.1 shows the power as function of α for the exact Monte-Carlo p-value, the conservative p-value and $Q_T(0.5)/0.5$ for $\rho = 0.5$.

The power of the de-randomized p-values is considerably better than that of the conservative. The choice of β has not a big influence on the power, therefore the adaptive version brings no benefit.

4.3.2. K-by-L Tables

Now suppose that we have multiple K -by- L tables $S^{(i)}$ with $\max(K, L) > 2$ and we want to test for independence. To get an exact p-value for one table $S^{(i)}$, we compute Pearson's chi-squared statistic $T_s^{(i)}$ for all t tables with the same marginal totals as $S^{(i)}$. The p-value is given by

$$\pi_i^{\text{cons}} := \frac{\#\{s \leq t: T_s^{(i)} \geq T_0^{(i)}\}}{t},$$

where $T_0^{(i)}$ is the statistic of $S^{(i)}$. If t is small, π_i^{cons} can take only few values and therefore it would be worthwhile to consider the randomized p-value

$$\pi_i \sim \text{Unif} \left[\frac{\#\{s \leq t: T_s^{(i)} > T_0^{(i)}\}}{t}, \frac{\#\{s \leq t: T_s^{(i)} \geq T_0^{(i)}\}}{t} \right].$$

For multiple tables these p-values can be combined using (4.1). Alternatively, noting that $T_0^{(i)} \sim_{\text{appr.}} \chi_{(K-1)(L-1)}^2$, we define the combined p-value

$$\pi := 1 - F_{m(K-1)(L-1)} \left(\sum_{i=1}^m F_{(K-1)(L-1)}^{-1}(\pi_i) \right), \quad (4.3)$$

where F_k denotes the c.d.f. of χ_k^2 .

Example 4.6. Table 4.2 shows data from a hypothetical multi-center clinical trial. For each of the five 2-by-3 tables the p-value π_i^{cons} is given. Combining them with (4.1), we get $\pi^{\text{cons}} = 0.122$. If we combine the randomized p-values π_i using (4.1) and de-randomize the result, we end up with a considerably smaller p-value. For example $Q_T(0.1)/0.1 = 0.011$.

Combining the p-values with (4.3) leads to even better results, namely $\pi^{\text{cons}} = 0.051$ and $Q_T(0.1)/0.1 = 0.008$.

Table			p-value
0	3	3	0.3
3	0	0	
0	2	1	0.333
3	0	0	
0	0	3	0.5
2	1	1	
0	3	0	0.111
6	0	2	
0	0	2	0.333
3	4	1	

Table 4.2.: Data from a hypothetical multi-center clinical trial.

Appendix

A. Classical Results

A.1. Lindeberg-Feller Central Limit Theorem

Theorem A.1. For $n = 1, 2, 3, \dots$ let $\mathbf{Y}_{n,1}, \mathbf{Y}_{n,2} \dots \mathbf{Y}_{n,n} \in \mathbb{R}^d$ independent random vectors. Suppose that for a matrix $\Sigma \in \mathbb{R}^{d \times d}$

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}(\mathbf{Y}_{n,i}) &= \mathbf{0}, \\ \sum_{i=1}^n \text{Var}(\mathbf{Y}_{n,i}) &\rightarrow \Sigma, \\ \sum_{i=1}^n \mathbb{E}(\|\mathbf{Y}_{n,i}\|^2 \min(1, \|\mathbf{Y}_{n,i}\|)) &\rightarrow 0. \end{aligned}$$

Then $\sum_{i=1}^n \mathbf{Y}_{n,i} \rightarrow_{\mathcal{L}} \mathcal{N}_d(\mathbf{0}, \Sigma)$.

This is a standard result in asymptotic statistics, see e.g. van der Vaart (1998).

A.2. Neyman-Pearson Lemma

Theorem A.2. Let P_0 and P_1 be probability distributions possessing densities f_0 and f_1 , respectively, with respect to a measure μ .

- (i) Existence. For testing $H : P_0$ against the alternative $K : P_1$ there exists a test φ and constants $c \in [0, \infty]$ and $\gamma \in [0, 1]$ such that

$$\mathbb{E}_0 \varphi(X) = \alpha \tag{A.1}$$

and

$$\varphi(x) = \begin{cases} 1 & \text{if } f_1(x) > cf_0(x) \\ \gamma & \text{if } f_1(x) = cf_0(x) \\ 0 & \text{if } f_1(x) < cf_0(x). \end{cases} \tag{A.2}$$

- (ii) Sufficient condition for a most powerful test. If a test satisfies (A.1) and (A.2) for some c and γ , then it is most powerful for testing P_0 against P_1 at level α .

A. Classical Results

- (ii) Necessary condition for the most powerful test. *If φ is most powerful at level α for testing P_0 against P_1 , then for some c it satisfies*

$$\varphi(x) = \begin{cases} 1 & \text{if } f_1(x) > cf_0(x) \\ 0 & \text{if } f_1(x) < cf_0(x) \end{cases} \quad (\text{A.3})$$

a.s. μ . It also satisfies (A.1) unless there exists a test of size $< \alpha$ and with power 1.

The proof can be found in Shao (2003).

References

- Agresti, A. (2007). *An introduction to categorical data analysis*. Wiley Series in Probability and Statistics. Wiley-Interscience [John Wiley & Sons], Hoboken, NJ, second edition.
- Bache, K. and Lichman, M. (2013). UCI machine learning repository.
- Bates, D. and Maechler, M. (2010). *Matrix: Sparse and Dense Matrix Classes and Methods*. R package version 0.999375-44.
- Dümbgen, L. (1998). On Tyler’s M -functional of scatter in high dimension. *Ann. Inst. Statist. Math.*, 50(3):471–491.
- Dümbgen, L. (2010). *Empirische Prozesse*. Lecture Notes, University of Bern.
- Dümbgen, L. (2011). *Multivariate Statistik*. Lecture Notes, University of Bern.
- Dümbgen, L., Igl, B.-W., and Munk, A. (2008). P-values for classification. *Electron. J. Stat.*, 2:468–493.
- Dümbgen, L., Pauly, M., and Schweizer, T. (2013). A survey of m -functionals of multivariate location and scatter. Technical Report 77, University of Bern. arXiv:1312.5594.
- Friedman, J. H., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22.
- Kiefer, J. (1953). Sequential minimax search for a maximum. *Proceedings of the American Mathematical Society*, 4(3):pp. 502–506.
- Kushmerick, N. (1999). Learning to remove internet advertisements. In *Proceedings of the third annual conference on Autonomous Agents*, AGENTS ’99, pages 175–181, New York, NY, USA. ACM.
- Lehmann, E. L. (1959). *Testing statistical hypotheses*. John Wiley & Sons Inc., New York.

References

- McLachlan, G. J. (1992). *Discriminant analysis and statistical pattern recognition*. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics. John Wiley & Sons Inc., New York. A Wiley-Interscience Publication.
- Meinshausen, N., Meier, L., and Bühlmann, P. (2009). p -values for high-dimensional regression. *J. Amer. Statist. Assoc.*, 104(488):1671–1681.
- Muirhead, R. J. (1982). *Aspects of multivariate statistical theory*. John Wiley & Sons Inc., New York. Wiley Series in Probability and Mathematical Statistics.
- R Core Team (2014). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Shao, J. (2003). *Mathematical statistics*. Springer Texts in Statistics. Springer-Verlag, New York, second edition.
- Taylor, A. E. and Lay, D. C. (1980). *Introduction to functional analysis*. John Wiley & Sons, New York-Chichester-Brisbane, second edition.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1):267–288.
- van der Vaart, A. W. (1998). *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak convergence and empirical processes*. Springer Series in Statistics. Springer-Verlag, New York. With applications to statistics.
- Zhu, J. and Hastie, T. (2004). Classification of gene microarrays by penalized logistic regression. *Biostatistics*, 5(14):427–443.
- Zumbrunnen, N. (2009). P -values for weighted nearest-neighbor classifiers. Master’s thesis, University of Bern.
- Zumbrunnen, N. and Dümmbgen, L. (2011). *pvclass: P-values for Classification*. R package version 1.01.

Index

- central limit theorem
 - inclusion probabilities, 70
 - Lindeberg-Feller, 97
 - missclassification rates, 59, 60
- class label, 3, 15
- classification, 3
- classifier, 4
 - optimal, 4, 59
 - standard linear, 6, 59
- contingency tables
 - multiple K -by- L , 92
 - multiple two-by-two, 88
- data set
 - buerk, 21, 46
 - interned ad, 45
 - mushrooms, 45
- distribution
 - elliptically symmetric, 6
 - Gaussian, 4
 - multivariate t , 49
 - spherically symmetric, 6
- empirical measure, 7, 54
- empirical process, 54
- exchangeability, 15
- feature space, 3
- feature vector, 3
- Gaussian model, *see* standard model
- golden section search
 - extended, 40
- half-space, 48
- inclusion probabilities, 70
 - conditional, 17
 - conditional empirical, 17
- Kronecker product, 28
- linear discriminant analysis, 6
- logistic regression
 - penalized multi-category, 17, 27
- M -estimator, 6
 - symmetrized, 7
- Mahalanobis distance, 4, 19
- missclassification rates, 9
- multiple testing, 11
- multiple use, 11
- nearest neighbors, 16, 19
 - k , 7
 - weighted, 9
- Neyman-Pearson Lemma, 97
- p-value, 10
 - combined, 86
 - cross-validated, 17
 - Monte-Carlo, 90
 - nonparametric, 15, 16
 - optimal, 13, 14
 - randomized, 84
- pattern probabilities
 - conditional, 18
 - empirical, 17

Index

plug-in statistic, *see* standard
model

posterior distribution, 9

posterior weights, 9

prediction region, 10, 15

prior probability, 3

pvcclass, 20

ROC curves

empirical, 18, 26

shortcut, 21

single use, 11

stability, 39

standard estimators, 6

standard model, 4, 15

plug-in statistic, 16, 19

subsampling, 40

training data, 5

training data, 15

tuning parameter, 39

unimodal, 40

List of Symbols

Symbols used in Chapters 1–3. This list is not exhaustive.

$B(\mathbf{x}, r)$	closed ball of radius r centered at \mathbf{x} , p. 7
$D_{\Sigma}(\mathbf{x}, \mathbf{y})$	Mahalanobis distance between $\mathbf{x} \in \mathbb{R}$ and $\mathbf{y} \in \mathbb{R}$ with respect to Σ , p. 4
F_{ξ}	distribution function of the one-dimensional random variable ξ , p. 48
$H(\beta, \gamma)$	half-space in \mathbb{R}^d , p. 48
$I(\cdot)$	ordered elements of \mathcal{G}_{θ} , p. 15
L	number of classes, p. 3
M	measure on \mathcal{X} , p. 3, 15
N_{θ}	number of training observations of class θ , p. 5, 15
P_{θ}	conditional distribution $\mathcal{L}(\mathbf{X} \mid Y = \theta)$, p. 3, 15
R_{θ}	missclassification rate, p. 9
R	risk of missclassification, p. 4
$S(\tau, \mathbf{X}, \theta)$	sum of test statistics, p. 39
$T_{\theta}(\mathbf{X}, \mathcal{D})$	test statistic based on training data, p. 16
$T_{\theta}^*(\mathbf{x})$	test statistic for the optimal p-value, p. 14, 15
$U(\mathbf{x}, r)$	open ball of radius r centered at \mathbf{x} , p. 7
$W_n(i)$	weight assigned to observation i , p. 9
Y	class label, p. 3, 15, 47
Z_1	first component of \mathbf{Z} , p. 47
$\mathbb{1}(A)$	indicator function for the set A , p. 4
\mathbf{A}	general matrix notation, $\mathbf{A} = \begin{pmatrix} A_{1,1} & \cdots & A_{1,d} \\ \vdots & \ddots & \vdots \\ A_{d,1} & \cdots & A_{d,d} \end{pmatrix}$, p. 48
$\Delta_{\mathbf{A}}$	scaled difference $\sqrt{n}(\hat{\mathbf{A}} - \mathbf{A})$, p. 48
$\Delta_{\mathbf{v}}$	scaled difference $\sqrt{n}(\hat{\mathbf{v}} - \mathbf{v}) = (\Delta_{\mathbf{v},1}, \dots, \Delta_{\mathbf{v},d})^{\top}$, p. 48
\mathbb{I}_d	d -dimensional identity matrix, p. 47
$\mathcal{N}_d(\boldsymbol{\mu}_{\theta}, \Sigma)$	d -dimensional Gaussian distribution with mean vector $\boldsymbol{\mu}_{\theta} \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$, p. 4
P_0	distribution of \mathbf{Z} , p. 47
Σ	positive definite covariance matrix in $\mathbb{R}^{d \times d}$, p. 4, 47
U^{\perp}	orthogonal complement of U , p. 47

\mathbf{X}	feature vector, p. 3, 15, 47
$\mathbf{Y}_{n,i}^{\mathbf{A}}$	summand of $\Delta_{\mathbf{A}}$, p. 48
$\mathbf{Y}_{n,i}^{\mathbf{v}}$	summand of $\Delta_{\mathbf{v}}$, p. 48
\mathbf{Z}	$\Sigma^{-1/2}(\mathbf{X}_1 - \boldsymbol{\mu}_1)$, p. 47
$\boldsymbol{\beta}$	$\Sigma^{-1/2}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)$, such that $\ \boldsymbol{\beta}\ = D_{\Sigma}(\boldsymbol{\mu}_2, \boldsymbol{\mu}_1)$; in the proof of Lemma 3.13 $\boldsymbol{\beta} := \Sigma^{-1/2}(\boldsymbol{\mu}_{\lambda} - \boldsymbol{\mu}_{\theta})$, p. 47, 73
$\mathbf{1}_d$	d -dimensional vector of ones, p. 31
\mathbf{e}_1	first standard unit vector, p. 50
$\ \cdot\ _{\infty}$	uniform norm, p. 55
$\ \cdot\ $	Euclidean or Frobenius norm for vectors or matrices, respectively, p. 5, 32
\hat{R}_{θ}	cross-validated missclassification rate, p. 9
$\hat{\mathbf{A}}$	estimator for \mathbf{A} , p. 48
$\hat{\Sigma}_M$	M -estimator, p. 6
$\hat{\Sigma}_{sym}$	symmetrized M -estimator, p. 6
$\hat{\Sigma}$	standard estimator for Σ or one of the M -estimators $\hat{\Sigma}_M$ and $\hat{\Sigma}_{sym}$, p. 6
$\hat{\boldsymbol{\mu}}_{\theta}$	standard estimator for $\boldsymbol{\mu}_{\theta}$, p. 6
$\hat{\mathbf{v}}$	estimator for \mathbf{v} , p. 48
\hat{P}	empirical measure of a sample of independent random variables with distribution P , p. 7, 54
$\hat{R}(\mathbf{x}, \mathbf{X}_i)$	rank of training observation \mathbf{X}_i , p. 9
$\hat{Y}(\mathbf{X})$	classifier, point predictor for Y , p. 4
$\hat{Y}^*(\mathbf{X})$	optimal classifier, p. 4
\hat{f}_{θ}	estimator of f_{θ} , p. 5
$\hat{r}_{k,n}(\mathbf{x}, \mathcal{D})$	radius of the smallest ball centered at \mathbf{x} , which covers at least k training vectors \mathbf{X}_i , p. 7
$\hat{w}_{\theta}(\mathbf{x}, \mathcal{D})$	estimator of $w_{\theta}(\mathbf{x})$, p. 16
\hat{w}_{θ}	estimator of w_{θ} , p. 5, 47
$\hat{\mathcal{I}}_{\alpha}(b, \theta)$	empirical conditional inclusion probability, p. 17, 70
$\hat{\mathcal{P}}_{\alpha}(b, S)$	empirical pattern probability, p. 17
$\hat{\mathcal{J}}_{\alpha}(\mathbf{X})$	prediction region, p. 10
$\hat{\mathcal{Y}}_{\alpha}(\mathbf{X}, \mathcal{D})$	prediction region based on training data, p. 15
$\mathbb{B}_{P,n}$	empirical process induced by \mathcal{H} , p. 54
$\boldsymbol{\mu}_{\theta,b}$	$(\boldsymbol{\mu}_{\theta} + \boldsymbol{\mu}_b)/2$, p. 5
$\boldsymbol{\mu}_{\theta}$	mean vector in \mathbb{R}^d , p. 4
\otimes	Kronecker product, p. 29, 48
$\pi_{\theta}(\mathbf{X}, \mathcal{D})$	p-value based on training data, p. 10, 16
$\pi_{\theta}(\mathbf{X}_i, \mathcal{D}_i)$	cross-validated p-value, p. 17

$\pi_\theta^*(\mathbf{x})$	optimal p-value, p. 13, 14
$\boldsymbol{\pi}$	vector of p-values $(\pi_\theta)_{\theta \in \mathcal{Y}}$, p. 13
\triangle	symmetric difference, p. 48
τ^*	optimal parameter τ , p. 40
$\tilde{\mathbf{X}}_i$	centered observation, p. 47
$\rightarrow_{\text{a.s.}}$	convergence in law, p. 48
\rightarrow_p	convergence in probability, p. 18, 48
$\rightarrow_{\mathcal{L}}$	convergence in law, p. 48
\mathbf{u}	$\ \boldsymbol{\beta}\ ^{-1}\boldsymbol{\beta}$, p. 59
$\mathbf{v}_{i:j}$	vector consisting of the components i to j of \mathbf{v} , p. 48
\mathbf{v}	general vector notation, $\mathbf{v} = (v_1, v_2, \dots, v_d)^\top$, p. 48
$\text{vec}(\mathbf{M})$	vector which is formed by stacking the columns of a matrix \mathbf{M} (from left to right), p. 34, 48
$d(\cdot, \cdot)$	some metric, p. 7
$f(\mathbf{x})$	density of the random vector \mathbf{X} , p. 11
f_θ	density of P_θ with respect to M , p. 3, 15
f_ξ	density of the random variable ξ , p. 48
g_θ	continuous bounded function such that $g_\theta(r)r^2$ is bounded for $r \geq 0$, p. 56
h_θ	continuous bounded function, p. 56
n	training sample size, p. 15, 47
$w_\theta(\mathbf{x})$	posterior weight, p. 9, 14
w_θ	prior probability, p. 3
$w_{b,\theta}$	ratio of prior weights, p. 14
$z_{1-\alpha/2}$	$(1 - \alpha)$ -quantile of the standard Gaussian distribution, p. 71
$\mathcal{D}(\mathbf{X}, \theta)$	training data extended by (\mathbf{X}, θ) , p. 21
$\mathcal{D}_i(\mathbf{X}_i, \mathbf{X}, \theta)$	training data after adding the observation (\mathbf{X}, θ) and with the class label of observation \mathbf{X}_i set to θ , p. 39
$\mathcal{D}_i(\mathbf{x})$	training data with \mathbf{x} in place of \mathbf{X}_i , p. 16
\mathcal{D}_i	training data without observation (\mathbf{X}_i, Y_i) , p. 9, 17
\mathcal{D}	training data, consisting of pairs (\mathbf{X}_i, Y_i) , for $i = 1, \dots, n$, p. 5
\mathcal{G}_θ	index set of training observations of class θ , p. 5
\mathcal{H}	collection of all half-spaces in \mathbb{R}^d , p. 48
$\mathcal{I}_\alpha(b, \theta \mid \mathcal{D})$	conditional inclusion probability, p. 17, 70
$\mathcal{L}(\mathbf{X})$	distribution of the random variable \mathbf{X} , p. 3
$\mathcal{P}_\alpha(b, S \mid \mathcal{D})$	pattern probability, p. 18
\mathcal{R}_α	measure of risk for p-values, p. 13
\mathcal{X}_0	support of $\mathcal{L}(\mathbf{X})$, p. 7
\mathcal{X}	feature space, p. 3, 15

List of Symbols

\mathcal{Y} set of class labels $\{1, \dots, L\}$, p. 3

Erklärung

gemäss Art. 28 Abs. 2 RSL 05

Name/Vorname: Zumbrunnen Niki

Matrikelnummer: 04-124-269

Studiengang: Statistik, Dissertation

Titel der Arbeit: P-Values for Classification – Computational Aspects and Asymptotics

Leiter der Arbeit: Prof. Dr. L. Dümbgen und Prof. Dr. A. Munk

Ich erkläre hiermit, dass ich diese Arbeit selbständig verfasst und keine anderen als die angegebenen Quellen benutzt habe. Alle Stellen, die wörtlich oder sinngemäss aus Quellen entnommen wurden, habe ich als solche gekennzeichnet. Mir ist bekannt, dass andernfalls der Senat gemäss Artikel 36 Absatz 1 Buchstabe r des Gesetzes vom 5. September 1996 über die Universität zum Entzug des auf Grund dieser Arbeit verliehenen Titels berechtigt ist.

Bern, 05.03.2014

Niki Zumbrunnen

Lebenslauf

1985	Geboren am 5. September in Bern
1992–1998	Primarschule in Bern
1998–2000	Sekundarschule in Bern
2000–2004	Mathematisch-Naturwissenschaftliches Gymnasium Bern-Neufeld (Schwerpunktfach Physik und angewandte Mathematik)
2004–2008	Bachelorstudium in Mathematik an der Universität Bern Minors: Philosophie und Naturwissenschaften Computational Science Informatik
2006–2009	Masterstudium in Mathematik an der Universität Bern
2009–2014	Doktorat an den Universitäten Bern und Göttingen